

Table of Contents

Index.....	1
Part 1 - Introduction - The Software Stack.....	4
Part 2 - GPU Memory Architecture and the Command Processor.....	12
Part 3 - 3D Pipeline Overview, Vertex Processing.....	18
Part 4 - Texture Samplers.....	24
Part 5 - Primitive Assembly, Clip and Cull, Projection, and Viewport Transform.....	34
Part 6 - (Triangle) Rasterization and Setup.....	42
Part 7 - Z Stencil Processing, 3 Different Ways.....	48
Part 8 - Pixel Processing (Fork Phase).	54
Part 9 - Pixel Processing (Join Phase).	62
Part 10 - Geometry Shaders.....	68
Part 11 - Stream-Out.....	72
Part 12 - Tessellation.....	74
Part 13 - Compute Shaders.....	81

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011: Index

July 9, 2011

Welcome.

This is the index page for a series of blog posts I'm currently writing about the D3D/OpenGL graphics pipelines *as actually implemented by GPUs*. A lot of this is well known among graphics programmers, and there's tons of papers on various bits and pieces of it, but one bit I've been annoyed with is that while there's both broad overviews and very detailed information on individual components, there's not much in between, and what little there is is mostly out of date.

This series is intended for graphics programmers that know a modern 3D API (at least OpenGL 2.0+ or D3D9+) well and want to know how it all looks under the hood. It's *not* a description of the graphics pipeline for novices; if you haven't used a 3D API, most if not all of this will be completely useless to you. I'm also assuming a working understanding of contemporary hardware design – you should at the very least know what registers, FIFOs, caches and pipelines are, and understand how they work. Finally, you need a working understanding of at least basic parallel programming mechanisms. A GPU is a massively parallel computer, there's no way around it.

Some readers have commented that this is a really low-level description of the graphics pipeline and GPUs; well, it all depends on where you're standing. GPU architects would call this a *high-level* description of a GPU. Not quite as high-level as the multicolored flowcharts you tend to see on hardware review sites whenever a new GPU generation arrives; but, to be honest, that kind of reporting tends to have a very low information density, even when it's done well. Ultimately, it's not meant to explain how anything actually *works* – it's just technology porn that's trying to show off shiny new gizmos. Well, I try to be a bit more substantial here, which unfortunately means less colors and less benchmark results, but instead lots and lots of text, a few monocolored diagrams and even some (*shudder*) equations. If that's okay with you, then here's the index:

Part 1 (<https://fgiesen.wordpress.com/2011/07/01/a-trip-through-the-graphics-pipeline-2011-part-1/>): Introduction; the Software stack.

Part 2 (<https://fgiesen.wordpress.com/2011/07/02/a-trip-through-the-graphics-pipeline-2011-part-2/>): GPU memory architecture and the Command Processor.

Part 3 (<https://fgiesen.wordpress.com/2011/07/03/a-trip-through-the-graphics-pipeline-2011-part-3/>): 3D pipeline overview, vertex processing.

Part 4 (<https://fgiesen.wordpress.com/2011/07/04/a-trip-through-the-graphics-pipeline-2011-part-4/>): Texture samplers.

Part 5 (<https://fgiesen.wordpress.com/2011/07/05/a-trip-through-the-graphics-pipeline-2011-part-5/>): Primitive Assembly, Clip/Cull, Projection, and Viewport transform.

Part 6 (<https://fgiesen.wordpress.com/2011/07/06/a-trip-through-the-graphics-pipeline-2011-part-6/>): (Triangle) rasterization and setup.

Part 7 (<https://fgiesen.wordpress.com/2011/07/08/a-trip-through-the-graphics-pipeline-2011-part-7/>): Z/Stencil processing, 3 different ways.

Part 8 (<https://fgiesen.wordpress.com/2011/07/10/a-trip-through-the-graphics-pipeline-2011-part-8/>): Pixel processing – “fork phase”.

Part 9 (<https://fgiesen.wordpress.com/2011/07/12/a-trip-through-the-graphics-pipeline-2011-part-9/>): Pixel processing – “join phase”.

Part 10 (<https://fgiesen.wordpress.com/2011/07/20/a-trip-through-the-graphics-pipeline-2011-part-10/>): Geometry Shaders.

Part 11 (<https://fgiesen.wordpress.com/2011/08/14/a-trip-through-the-graphics-pipeline-2011-part-11/>): Stream-Out.

Part 12 (<https://fgiesen.wordpress.com/2011/09/06/a-trip-through-the-graphics-pipeline-2011-part-12/>): Tessellation.

Part 13 (<https://fgiesen.wordpress.com/2011/10/09/a-trip-through-the-graphics-pipeline-2011-part-13/>): Compute Shaders.



(<http://creativecommons.org/publicdomain/zero/1.0/>)

To the extent possible under law,

Fabian Giesen (<https://fgiesen.wordpress.com>)

has waived all copyright and related or neighboring rights to

A trip through the Graphics Pipeline 2011.

From → Coding, Graphics Pipeline

38 Comments**1. nordicsavage permalink**

Don't forget a chapter on Multi-sampling and the difference between the various AA techniques out there.. :) Am loving this series though :)

Reply

- o **fgiesen permalink**

Unlikely. I'm sticking with the basic D3D11 pipeline, and even there I'm dropping some subjects. Basic MSAA (2x, 4x, 8x) must be supported by every D3D11 device, all fancier stuff is strictly optional. And there's even bits in core D3D11 that I'm mostly ignoring – point and line primitives (and their setup+rasterization), the finer points of cube map filtering, the various trilinear filtering "optimizations" that GPUs do (not, ahem, strictly in accordance with the spec...), and so on...

All that's there, but I have enough material to write another 7 parts already; I do plan to finish this series eventually :)

Reply

2. Manu permalink

Good job man!

Maybe you should write a book or something. Many people will buy it for sure.

Reply

- o **doctor_shim permalink**

that would negatively impact the size of the readership, in addition to reducing exposure!

Reply

- o **fgiesen permalink**

Cleaned up, expanded PDF version is in the works. You're gonna have to print it yourself though. :)

3. Brandon Furtwangler permalink

This is a great series of articles. Thanks for making them. Can't wait for the part(s) on compute shaders.

I'd love to hear how you think the pipeline could/should evolve in the coming years.

Reply

4. Francis Boivin permalink

Are you considering talking about constant buffers? They are an important part of d3d10+. At least, they are from an API perspective – do driver actually still care about this optimized representation of shader parameters? I'm not an OpenGL guy so I don't know if it went with a similar API or if constants still set using glUniform[...]?

Reply

- o **fgiesen permalink**

They're an important part of the API, but on the GPU side they're really just chunks of memory that the shader units can access.

Originally, the constants used to be an actual special register file on the chip. D3D10 increased the limits too far for that to be practical: up to 16 constant buffers per stage – 15 API-accessible and one for immediate constants – with up to 4k elements of up to 16 bytes each; so up to 64k per CB and up to 1MB total, *per stage* – and you can have a lot of them active at once (worst case, five of them: VS, HS, DS, GS, PS). So now CBs are just regular buffers in GPU/host memory (like VBs, IBs or Textures) and can be accessed as such. One option is to still have some (smaller) amount of fast memory reserved for constants, and try to pack the CBs that will fit in there. But this generation of GPUs has (this is fairly new!) a regular fast cache between shader units and memory. With a cache, I'd just leave the CBs completely memory-mapped and let the cache deal with it! This can adapt to the dynamic behavior of shaders, rather than having to rely on some heuristic to pick which CBs to pack into fast memory.

Reply

5. ikrima permalink

Dude, this is amazing. I took a sabbatical out of gfx/vfx for 4 years and been playing catch up over the last 7 months. Thanks for the awesome in-depth articles for an easy quick dive into how things have changed.

Reply

6. sinistraldexter permalink

Reblogged this on sinistraldexter and commented:
one of the best work w.r.t the graphics pipeline

Reply

7. nandu permalink

Hi, intern@ Nvidia. Really helpful to get started.

Reply

8. Samuel Egger permalink

An amazing series! Thank you. Although I am quiet curious where one learns all these things?

Reply

- o **fgiesen permalink**

I was working on a GPU at the time I wrote the series.

Reply

9. Suso permalink

I do you plan to review the new mesh and amplification shaders?

Reply

Trackbacks & Pingbacks

1. Real-Time Rendering · Seven Things for July 24th, 2011
2. Viaje alucinante por un pipeline grafico « martin b.r.
3. A trip through the Graphics Pipeline | Light is beautiful
4. A Very Good Technical Guide to the 3D Graphics Pipeline
5. A trip through the Graphics Pipeline 2011, part 1 « The ryg blog
6. A trip through the Graphics Pipeline 2011, part 5 « The ryg blog
7. A trip through the Graphics Pipeline 2011, part 4 « The ryg blog
8. A trip through the Graphics Pipeline 2011, part 6 « The ryg blog
9. A trip through the Graphics Pipeline 2011, part 7 « The ryg blog
10. A trip through the Graphics Pipeline 2011, part 8 « The ryg blog
11. A trip through the Graphics Pipeline 2011, part 3 « The ryg blog
12. A trip through the Graphics Pipeline 2011, part 2 « The ryg blog
13. A trip through the Graphics Pipeline 2011, part 9 « The ryg blog
14. A trip through the Graphics Pipeline 2011, part 10 « The ryg blog
15. A trip through the Graphics Pipeline 2011, part 11 « The ryg blog
16. A trip through the Graphics Pipeline 2011, part 13 « The ryg blog
17. A trip through the Graphics Pipeline 2011, part 12 « The ryg blog
18. Programming | PearlTrees
19. Túra a grafikus csővezetékben | cikksorozat | szimpatikus.hu trackback proxy
20. Xbox / PC, early-Z and early stencil in XNA « IceFall Games
21. Order and types of depth testing « Interplay of Light
22. HPG 2013 | dickyjim
23. A trip through the Graphics Pipeline | The blog at the bottom of the sea
24. What's the big deal with Apples Metal API? | RenderingPipeline

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 1

July 1, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

It's been awhile since I posted something here, and I figured I might use this spot to explain some general points about graphics hardware and software as of 2011; you can find functional descriptions of what the graphics stack in your PC does, but usually not the “how” or “why”; I'll try to fill in the blanks without getting too specific about any particular piece of hardware. I'm going to be mostly talking about DX11-class hardware running D3D9/10/11 on Windows, because that happens to be the (PC) stack I'm most familiar with – not that the API details etc. will matter much past this first part; once we're actually on the GPU it's all native commands.

The application

This is your code. These are also your bugs. Really. Yes, the API runtime and the driver have bugs, but this is not one of them. Now go fix it already.

The API runtime

You make your resource creation / state setting / draw calls to the API. The API runtime keeps track of the current state your app has set, validates parameters and does other error and consistency checking, manages user-visible resources, may or may not validate shader code and shader linkage (or at least D3D does, in OpenGL this is handled at the driver level) maybe batches work some more, and then hands it all over to the graphics driver – more precisely, the user-mode driver.

The user-mode graphics driver (or UMD)

This is where most of the “magic” on the CPU side happens. If your app crashes because of some API call you did, it will usually be in here :). It's called “nvd3dum.dll” (Nvidia) or “atiumd*.dll” (AMD). As the name suggests, this is user-mode code; it's running in the same context and address space as your app (and the API runtime) and has no elevated privileges whatsoever. It implements a lower-level API (the DDI) that is called by D3D; this API is fairly similar to the one you're seeing on the surface, but a bit more explicit about things like memory management and such.

This module is where things like shader compilation happen. D3D passes a pre-validated shader token stream to the UMD – i.e. it's already checked that the code is valid in the sense of being syntactically correct and obeying D3D constraints (using the right types, not using more textures/samplers than available, not exceeding the number of available constant buffers, stuff like that). This is compiled from HLSL code and usually has quite a number of high-level optimizations (various loop optimizations, dead-code elimination, constant propagation, predication ifs etc.) applied to it – this is good news since it means the driver benefits from all these relatively costly optimizations that have been performed at compile time. However, it also has a bunch of lower-level optimizations (such as register allocation and loop unrolling) applied that drivers would rather do themselves; long story short, this usually just gets immediately turned into a intermediate representation (IR) and then compiled some more; shader hardware is close enough to D3D bytecode that compilation doesn't need to work wonders to give good results (and the HLSL compiler having done some of the high-yield and high-cost optimizations already definitely helps), but there's still lots of low-level details (such as HW resource limits and scheduling constraints) that D3D neither knows nor cares about, so this is not a trivial process.

And of course, if your app is a well-known game, programmers at NV/AMD have probably looked at your shaders and wrote hand-optimized replacements for their hardware – though they better produce the same results lest there be a scandal :). These shaders get detected and substituted by the UMD too. You’re welcome.

More fun: Some of the API state may actually end up being compiled into the shader – to give an example, relatively exotic (or at least infrequently used) features such as texture borders are probably not implemented in the texture sampler, but emulated with extra code in the shader (or just not supported at all). This means that there’s sometimes multiple versions of the same shader floating around, for different combinations of API states.

Incidentally, this is also the reason why you’ll often see a delay the first time you use a new shader or resource; a lot of the creation/compilation work is deferred by the driver and only executed when it’s actually necessary (you wouldn’t believe how much unused crap some apps create!). Graphics programmers know the other side of the story – if you want to make sure something is actually created (as opposed to just having memory reserved), you need to issue a dummy draw call that uses it to “warm it up”. Ugly and annoying, but this has been the case since I first started using 3D hardware in 1999 – meaning, it’s pretty much a fact of life by this point, so get used to it. :)

Anyway, moving on. The UMD also gets to deal with fun stuff like all the D3D9 “legacy” shader versions and the fixed function pipeline – yes, all of that will get faithfully passed through by D3D. The 3.0 shader profile ain’t that bad (it’s quite reasonable in fact), but 2.0 is crusty and the various 1.x shader versions are seriously whack – remember 1.3 pixel shaders? Or, for that matter, the fixed-function vertex pipeline with vertex lighting and such? Yeah, support for all that’s still there in D3D and the guts of every modern graphics driver, though of course they just translate it to newer shader versions by now (and have been doing so for quite some time).

Then there’s things like memory management. The UMD will get things like texture creation commands and need to provide space for them. Actually, the UMD just suballocates some larger memory blocks it gets from the KMD (kernel-mode driver); actually mapping and unmapping pages (and managing which part of video memory the UMD can see, and conversely which parts of system memory the GPU may access) is a kernel-mode privilege and can’t be done by the UMD.

But the UMD can do things like **swizzling textures** (<https://fgiesen.wordpress.com/2011/01/17/texture-tiling-and-swizzling/>) (unless the GPU can do this in hardware, usually using 2D blitting units not the real 3D pipeline) and schedule transfers between system memory and (mapped) video memory and the like. Most importantly, it can also write command buffers (or “DMA buffers” – I’ll be using these two names interchangeably) once the KMD has allocated them and handed them over. A command buffer contains, well, commands :). All your state-changing and drawing operations will be converted by the UMD into commands that the hardware understands. As will a lot of things you don’t trigger manually – such as uploading textures and shaders to video memory.

In general, drivers will try to put as much of the actual processing into the UMD as possible; the UMD is user-mode code, so anything that runs in it doesn’t need any costly kernel-mode transitions, it can freely allocate memory, farm work out to multiple threads, and so on – it’s just a regular DLL (even though it’s loaded by the API, not directly by your app). This has advantages for driver development too – if the UMD crashes, the app crashes with it, but not the whole system; it can just be replaced while the system is running (it’s just a DLL!); it can be debugged with a regular debugger; and so on. So it’s not only efficient, it’s also convenient.

But there’s a big elephant in the room that I haven’t mentioned yet.

Did I say “user-mode driver”? I meant “user-mode drivers”.

As said, the UMD is just a DLL. Okay, one that happens to have the blessing of D3D and a direct pipe to the KMD, but it’s still a regular DLL, and it runs in the address space of its calling process.

But we’re using multi-tasking OSes nowadays. In fact, we have been for some time.

This “GPU” thing I keep talking about? That’s a shared resource. There’s only one that drives your main display (even if you use SLI/Crossfire). Yet we have multiple apps that try to access it (and pretend they’re the only ones doing it). This doesn’t just work automatically; back in The Olden Days, the solution was to only give 3D to one app at a time, and while that app was active, all others wouldn’t have access. But that doesn’t really cut it if you’re trying to have your windowing system use the GPU for rendering. Which is why you need some component that arbitrates access to the GPU and allocates time-slices and such.

Enter the scheduler.

This is a system component – note the “the” is somewhat misleading; I’m talking about the graphics scheduler here, not the CPU or IO schedulers. This does exactly what you think it does – it arbitrates access to the 3D pipeline by time-slicing it between different apps that want to use it. A context switch incurs, at the very least, some state switching on the GPU (which generates extra commands for the command buffer) and possibly also swapping some resources in and out of video memory. And of course only one process gets to actually submit commands to the 3D pipe at any given time.

You’ll often find console programmers complaining about the fairly high-level, hands-off nature of PC 3D APIs, and the performance cost this incurs. But the thing is that 3D APIs/drivers on PC really have a more complex problem to solve than console games – they really do need to keep track of the full current state for example, since someone may pull the metaphorical rug from under them at any moment! They also work around broken apps and try to fix performance problems behind their backs; this is a rather annoying practice that no-one’s happy with, certainly including the driver authors themselves, but the fact is that the business perspective wins here; people expect stuff that runs to continue running (and doing so smoothly). You just won’t win any friends by yelling “BUT IT’S WRONG!” at the app and then sulking and going through an ultra-slow path.

Anyway, on with the pipeline. Next stop: Kernel mode!

The kernel-mode driver (KMD)

This is the part that actually deals with the hardware. There may be multiple UMD instances running at any one time, but there’s only ever one KMD, and if that crashes, then boom you’re dead – used to be “blue screen” dead, but by now Windows actually knows how to kill a crashed driver and reload it (progress!). As long as it happens to be just a crash and not some kernel memory corruption at least – if that happens, all bets are off.

The KMD deals with all the things that are just there once. There’s only one GPU memory, even though there’s multiple apps fighting over it. Someone needs to call the shots and actually allocate (and map) physical memory. Similarly, someone must initialize the GPU at startup, set display modes (and get mode information from displays), manage the hardware mouse cursor (yes, there’s HW handling for this, and yes, you really only get one! :), program the HW watchdog timer so the GPU gets reset if it stays unresponsive for a certain time, respond to interrupts, and so on. This is what the KMD does.

There’s also this whole content protection/DRM bit about setting up a protected/DRM’ed path between a video player and the GPU so no the actual precious decoded video pixels aren’t visible to any dirty user-mode code that might do awful forbidden things like dump them to disk (...whatever). The KMD has some involvement in that too.

Most importantly for us, the KMD manages the *actual* command buffer. You know, the one that the hardware actually consumes. The command buffers that the UMD produces aren’t the real deal – as a matter of fact, they’re just random slices of GPU-addressable memory. What actually happens with them is that the UMD finishes them, submits them to the scheduler, which then waits until that process is up and then passes the UMD command buffer on to the KMD. The KMD then writes a call to command buffer into the main command buffer, and depending on whether the GPU command processor can read from main memory or not, it may also need to DMA it to video memory first. The main command buffer is usually a (quite small) **ring buffer** (<https://fgiesen.wordpress.com/2010/12/14/ring-buffers-and-queues/>) – the only thing that ever gets written there is system/initialization commands and calls to the “real”, meaty 3D command buffers.

But this is still just a buffer in memory right now. Its position is known to the graphics card – there’s usually a read pointer, which is where the GPU is in the main command buffer, and a write pointer, which is how far the KMD has written the buffer yet (or more precisely, how far it has *told* the GPU it has written yet). These are hardware registers, and they are memory-mapped – the KMD updates them periodically (usually whenever it submits a new chunk of work)...

The bus

...but of course that write doesn’t go directly to the graphics card (at least unless it’s integrated on the CPU die!), since it needs to go through the bus first – usually PCI Express these days. DMA transfers etc. take the same route. This doesn’t take very long, but it’s yet another stage in our journey. Until finally...

The command processor!

This is the frontend of the GPU – the part that actually reads the commands the KMD writes. I'll continue from here in the next installment, since this post is long enough already :)

Small aside: OpenGL

OpenGL is fairly similar to what I just described, except there's not as sharp a distinction between the API and UMD layer. And unlike D3D, the (GLSL) shader compilation is not handled by the API at all, it's all done by the driver. An unfortunate side effect is that there are as many GLSL frontends as there are 3D hardware vendors, all of them basically implementing the same spec, but with their own bugs and idiosyncrasies. Not fun. And it also means that the drivers have to do all the optimizations themselves whenever they get to see the shaders – including expensive optimizations. The D3D bytecode format is really a cleaner solution for this problem – there's only one compiler (so no slightly incompatible dialects between different vendors!) and it allows for some costlier data-flow analysis than you would normally do.

Omissions and simplifications

This is just an overview; there's tons of subtleties that I glossed over. For example, there's not just one scheduler, there's multiple implementations (the driver can choose); there's the whole issue of how synchronization between CPU and GPU is handled that I didn't explain at all so far. And so on. And I might have forgotten something important – if so, please tell me and I'll fix it! But now, bye and hopefully see you next time.

From → Coding, Graphics Pipeline

39 Comments

1. Bitouo permalink

Thank you for this great article! Will you write a little bit about how synchronization between CPU and GPU is handled. I have been curious about it for a long while. Or maybe point out some nice articles I can read. :)

Reply

o fgiesen permalink

Yes, I'll be filling in the blanks about the details of memory/resource lifetime management and CPU/GPU synchronization as I go along. This will be one of our recurring themes in fact, since it affects virtually all parts of the pipeline in some way. I'll explain it as soon as I get there :)

Reply

2. Compulsive Dabbler permalink

Thanks a ton for this, I haven't found a single other resource that explains these details so concisely!

Reply

3. 3dfx permalink

Impressive article. Can't wait for the next installment!

Reply

4. Christophe Riccio permalink

"And it also means that the drivers have to do all the optimizations themselves whenever they get to see the shaders – including expensive optimizations."

This is actually a drawback of the D3D approach because each GPU architecture is really different, the first task of the compiler is to "un-optimized" the bytecode before running the GPU optimizations.

Reply

o fgiesen permalink

I'm actually working on a shader compiler for actual hardware right now (albeit a R&D one, not one intended for production) and I can tell you that – at least as far as my experience goes – this is an urban myth, or at least blown way out of proportion.

Yes, D3D bytecode is an abstraction, and I'd rather not have the HLSL compiler do "information-destroying" transforms like loop unrolling, function inlining or deciding whether to use ifs or predication. Detecting and undoing this does mean extra work for the compiler and is annoying and a waste of time, but in the grand scheme of things it's not that big a deal.

But the bulk of it is simply that the HLSL compiler does some optimizations which are "information-neutral" (such as register packing or scheduling) that are just a waste of time on the HLSL compiler side; if I compile a shader, the first thing I do is convert it to some optimizer-friendly IR (nowadays, that usually means Static Single Assignment form or its descendants), and that simply implicitly destroys these optimizations. There's no extra work involved in "un-optimizing" them – it just happens as a side effect of IR construction. So it's just a waste of time for HLSL to be doing this, but not actually a problem.

And it's definitely nice to have a shared frontend that does all the cool "code cleanup" transformations for me; the HLSL compiler does constant propagation, global value numbering, partial redundancy elimination and loop-invariant code motion, and has fairly sophisticated mechanisms for algebraic simplifications.

All of these things are costly both in terms of implementation complexity and run time, and they're more or less completely device independent. Any regular compiler does them before the IR reaches the device-dependent backend code generator (which is what the driver is); it makes perfect sense for D3D to be doing it too. It pulls expensive work out of the runtime compilation loop, and it's a lot of tricky work that drivers now don't need to worry about. In the OpenGL/GLSL world, they do, and if you've used GLSL particularly in its first few years, you got to see the (awful) results. I think farming this off into a separate stage was a very good call on MS's part.

Reply

5. KeyJ permalink

Great article. I have one question though: Where is the scheduler located? Is it part of the UMD, of D3D, the KMD or is it some completely different library or process?

Reply

- o **fgiesen permalink**

The scheduler is part of the OS/driver model.

The total sequence of stages up to this point is App (you) -> API (OS/driver) -> UMD (driver) -> Scheduler (OS) -> KMD (driver) -> GPU.

This also means that the notion of command/DMA buffer is part of the driver model, not just some implementation detail, since the scheduler needs to know about it (after all, that's the things it's scheduling!). It doesn't touch the data (or know what it means), but it does need to know that such a thing exists to pass it around.

Reply

- o **Marek permalink**

Is this sequence of stages applies to Linux also? I don't understand UMD idea in Linux.

- o **fgiesen permalink**

The structure is similar but not quite the same. The equivalent to the KMD is the kernel DRM driver, UMD is the GL driver, state tracker and everything up to and including libdrm/libdri. The details are different though. In particular, I'm not aware of a kernel-land central video memory manager, graphics scheduler or video memory paging / command buffer patching mechanism.

- o **x4da permalink**

In radeon queue scheduling and submission are done on kernel driver side. Same with memory management: userspace request kernel driver to set up VM address space and requests a Buffer Object to be created and mapped to VM AS, which kernel driver (using TTM mechanism) does and adds to its private BO pool.

- o **fgiesen permalink**

That's just the point though, this is part of the per-device drivers, not a central graphics subsystem service like on Windows.

- o **x4da permalink**

Yep, on linux thing could be different depending on vendor: proprietary nvidia drivers don't use kernel TTM buffer management and DRI/DRM interface.

6. Corbin Simpson permalink

A good overview.

Two notes: First, GL drivers generally are built on very reusable code these days, which means that the general pattern of GLSL->IR->driver-specific bytecode is actually very viable and happens in production drivers. Gallium takes this to an extreme: GLSL is compiled to a GLSL IR, which is optimized at the high level and passed to Gallium as TGSI, which is GPU-neutral and serializeable. Then the specific driver backend turns TGSI into actual GPU bytecode. The entire process works pretty well, actually!

Second, I figured I'd explain what's different on Linux, as far as the kernel goes. Not much changes; the big things are in what's shared and how robust the drivers are. GPU scheduling's per-driver and usually very basic; the predominant system involves lockless blocking ioctl(s) and a master process (X server, Wayland, etc.) which controls who is allowed to submit commands. On-GPU memory management is shared through a couple of kernel libraries, GEM and TTM, which allow everybody to enjoy common benefits in improvements to memory management. Finally, drivers can't be auto-evicted if they crash, but they do generally understand (at least for modern GPUs) how to reboot the GPU if it wedges or hangs.

Reply

- o [fgiesen permalink](#)

"the general pattern of GLSL->IR->driver-specific bytecode is actually very viable and happens in production drivers"
Yes, certainly; this general architecture is really the only sane way to build optimizing compilers when your backend architecture goes through significant changes every 1-2 years :). But D3D standardizing on a frontend/IR (even if it's not ideal) helped them get significantly less problems in "portability" of shader code across different vendors at a critical point in time.

Interesting to see that Linux is doing something very similar these days. Good stuff!

Reply

- o [przemo_li permalink](#)

Gallium 3D is idea of common front put even further. You can get others state trackers (eg. DX10, OpenVG, OpenGL ES, etc.) too. And it is evolving even more (eg. some devs want to use LLVM in the middle between GLSL and specific bytecode). Current DX model can not beat it. Since all common optimizations for hwd, must be implemented separately anyway. Gallium 3D is more of template you driver can fill with really unique stuff, while HLSL compiler is just top common stuff with some cleaning. (However, still its better than any non-gallium3d driver, but do not know how Mac OSX handle GLSL, maybe there Apple will copy HLSL approach?)

- o [fgiesen permalink](#)

As I've said in the article, the main reason I'm using the Win32/D3D graphics stack for this series (which is about GPUs not 3D APIs) is that it's the one I'm most familiar with, and I need to decide on *some* set of terminology to use. Beyond that, I'm not interested in talking about API differences. If you are, good for you, but please take it somewhere else.

7. Nico [permalink](#)

And don't forget about the additional cost of having the cpu flushing it's cache and draining the write buffers to memory first and the PCIe-Hostcontroller which has to read it back from there. This also increases latency.

Reply

- o [fgiesen permalink](#)

Command buffers are customarily in write-combined memory which is uncached (but store buffered). So no cache-draining involved; on x86, it's just a `sfence` (which flushes the store buffers). This isn't free, but considering you went through the OS scheduler and a user->kernel mode transition just to get the command buffers to the KMD in the first place, this is small fry. And actually, all of this is pipelined to the point where graphics drivers are *up to 3 frames ahead* of the GPU (drivers used to try for more, but that was totally screwing with input latency and such, so it's now capped at 3). So at that point (if the CPU/driver is fast enough to get that far ahead!), assuming 60Hz full-framerate rendering, the latency between an app issuing a draw call and the GPU actually processing is about 50ms, and can be south of that if we're not rendering at full framerate. Point being, there's so much intentional latency introduced by the SW stack just to keep the GPU fed at all times, that additional (comparatively small) latencies added by such low-level effects are all but unnoticeable.

Reply

8. Raja [permalink](#)

This is absolutely fantastic. Just what I needed! I've always struggled to see the big picture of interactions, and your blog is just what the doctor ordered.

Thanks a lot for taking the time to do this. I can't wait to read all your posts now!

Reply

9. dca [permalink](#)

A quick add: modern GPUs have somewhat called GPUVM – separate process address spaces with their own sets of page tables each.

Reply

10. HY permalink

What is the mechanism in which the UMD communicates with the KMD? Does the KMD create special syscalls which the UMD invokes? Or exposes a \GLOBAL?? device node to be accessed by the UMD?

Reply

- o [fgiesen permalink](#)

The D3D runtime provides an API that UMDs are using to talk to the KMD. Just refer to the official docs if you're interested in the details. I'm not sure how the implementation looks under the hood – graphics driver authors (both KMD and UMD) don't need to worry, this is mediated by the D3D runtime and the DirectX Graphics Kernel Subsystem (dxgkrnl).

Reply

11. James Bedford permalink

Great article – can't wait to keep reading more!

Reply

12. Yogesh permalink

Thanks for the great pipeline series. You should add one article on locks.(locking cpu or gpu memory).

Reply

13. Steex permalink

Thanks for the excellent series! I'm not a professional graphics programmer, still as a game developer I have to know at least the facade of the pipeline. But only the facade. So naturally while reading your series I was surprised... well, many times. Amazing algorithms! And pretty good explanation, I have to add.

It's interesting though, what is changed to the modern time (end of 2015). Are the articles still relevant? For example, I read that AMD Mantle and DirectX 12 provide direct access to command buffers. Did this change something in GPU architecture? Sorry for somewhat newbie questions. :)

Reply

- o [fgiesen permalink](#)

Neither Mantle, D3D12, nor Vulkan provide direct access to command buffers. They remove some of the intermediate layers in the software stack (by pushing the work to the application side), but that all happens before the stuff I talk about in this series.

There have been no fundamental changes to desktop GPU architecture since I wrote this series, not at the (still relatively abstract) level of this series anyway. A lot of the details have changed – for example, asynchronous compute support means that the GPU can process multiple command streams at the same time (by having multiple command processors, time-slicing a single command processor in either SW or HW, or some combination thereof). That means that from a user point of view, there's now multiple command processors; but that doesn't change the way they work, it just means there's more of them (complicating internal synchronization, but again, that's below the level of abstraction of this series). Another big-ticket item would be what D3D12 calls "Rasterizer Ordered Views"; but I already talked about several ways to handle blend ordering. Basically, ROV support means the GPU can optionally track in-flight quads that are trying to write to the same location, and make sure they run in order.

But for the most part, D3D12/Mantle/Vulkan are not about any of these things. The biggest change in these APIs is that they replace the "state machine" model of GL and older D3D versions with a model where all kinds of state is pushed by the application beforehand and compiled into a format the hardware can understand directly, and resource residency/memory/dependency management is pushed (to a significant extent) to the app. This saves a lot of work on the driver side, which otherwise has to do the state translation every time somebody changes any piece of state, and has to track which resources are in use by whom at what time to make sure the correct synchronization happens. That's where the speed-ups come from. The underlying HW didn't change at all.

Reply

14. cyang permalink

Hello fgiesen, thank you very much for this great article! I am a Ph.D student in EE and have been digging into Graphic cards as a side hobby project. Recently, I am thinking about making an open source PC graphic card project. Though this is a tough and unrealistic side project, it is still a good goal to move towards and learn stuff during the process. With that being said, I would appreciate it if you could allow me to translate this article into another language (in my case it is Chinese) and share these interesting series of article with some of my friends who are working with me on the side hobby project. Thank you very much and I look forward to hearing your permission. Have a nice day:)

Reply

- o [fgiesen permalink](#)

The series as a whole is CC-0 licensed, you are allowed to do whatever you want with it!

Reply

- **cyang permalink**

Thank you very much for your kind reply. Sorry I didn't check the series cover page before:)

15. Connor A. Haskins permalink

When you say the API runtime "manages user-visible resources", the user in this context is the application, right? What are examples of these resources?

Also, thank you so much for writing this.

Reply

- **fgiesen permalink**

Textures, buffers (and constant buffers), render targets are the major ones. "Resource" is D3D-speak (note I'm using D3D terminology the whole way through the series) for all of those.

Reply

Trackbacks & Pingbacks

1. Geeks3D Programming Links – July 01, 2011 - 3D Tech News, Pixel Hacking, Data Visualization and 3D Programming - Geeks3D.com
2. 3D Graphics Pipeline Explained - 3D Tech News, Pixel Hacking, Data Visualization and 3D Programming - Geeks3D.com
3. A trip through the Graphics Pipeline 2011: Index « The ryg blog
4. Understanding Modern GPUs (I): Introduction « TraxNet – Blog
5. Viaje alucinante por un pipeline grafico « martin b.r.
6. BreakTryCatch » Getting Started With DirectX 11
7. How do graphics processing units (GPUs) work? - Quora

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 2

July 2, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Not so fast.

In the previous part I explained the various stages that your 3D rendering commands go through on a PC before they actually get handed off to the GPU; short version: it’s more than you think. I then finished by name-dropping the command processor and how it actually finally does something with the command buffer we meticulously prepared. Well, how can I say this – I lied to you. We’ll indeed be meeting the command processor for the first time in this installment, but remember, all this command buffer stuff goes through memory – either system memory accessed via PCI Express, or local video memory. We’re going through the pipeline in order, so before we get to the command processor, let’s talk memory for a second.

The memory subsystem

GPUs don’t have your regular memory subsystem – it’s different from what you see in general-purpose CPUs or other hardware, because it’s designed for very different usage patterns. There’s two fundamental ways in which a GPU’s memory subsystem differs from what you see in a regular machine:

The first is that GPU memory subsystems are *fast*. Seriously fast. A Core i7 2600K will hit maybe 19 GB/s memory bandwidth – on a good day. With tail wind. Downhill. A GeForce GTX 480, on the other hand, has a total memory bandwidth of close to 180 GB/s – nearly an order of magnitude difference! Whoa.

The second is that GPU memory subsystems are *slow*. Seriously slow. A cache miss to main memory on a Nehalem (first-generation Core i7) takes about 140 cycles if you multiply the **memory latency as given by AnandTech** (<http://www.anandtech.com/show/2542/5>) by the clock rate. The GeForce GTX 480 I mentioned previously has a **memory access latency of 400-800 clocks** (<http://www.stanford.edu/dept/ICME/docs/seminars/Rennich-2011-04-25.pdf>). So let’s just say that, measured in cycles, the GeForce GTX 480 has a bit more than 4x the average memory latency of a Core i7. Except that Core i7 I just mentioned is clocked at 2.93GHz, whereas GTX 480 shader clock is 1.4 GHz – that’s it, another 2x right there. Woops – again, nearly an order of magnitude difference! Wait, something funny is going on here. My common sense is tingling. This must be one of those trade-offs I keep hearing about in the news!

Yep – GPUs get a massive increase in bandwidth, but they pay for it with a massive increase in latency (and, it turns out, a sizable hit in power draw too, but that’s beyond the scope of this article). This is part of a general pattern – GPUs are all about throughput over latency; don’t wait for results that aren’t there yet, do something else instead!

That’s almost all you need to know about GPU memory, except for one general DRAM tidbit that will be important later on: DRAM chips are organized as a 2D grid – both logically and physically. There’s (horizontal) row lines and (vertical) column lines. At each intersection between such lines is a transistor and a capacitor; if at this point you want to know how to actually build memory from these ingredients, **Wikipedia is your friend** (http://en.wikipedia.org/wiki/DRAM#Operation_principle). Anyway, the salient point here is that the address of a location in DRAM is split into a row address and a column address, and DRAM reads/writes internally always end up accessing all columns in the given row at the same time. What this means is that it’s much cheaper to access a swath of memory that maps to exactly one DRAM row than it is to access the same amount of memory spread across multiple rows. Right now this may seem like just a random bit of DRAM trivia, but this will become important later on; in other words, pay attention: this will be on the exam. But to tie this up with the figures in the previous paragraphs, just let me note that you can’t reach those peak memory bandwidth figures above by just reading a few bytes all over memory; if you want to saturate memory bandwidth, you better do it one full DRAM row at a time.

The PCIe host interface

From a graphics programmer standpoint, this piece of hardware isn't super-interesting. Actually, the same probably goes for a GPU hardware architect too. The thing is, you still start caring about it once it's so slow that it's a bottleneck. So what you do is get good people on it to do it properly, to make sure that doesn't happen. Other than that, well, this gives the CPU read/write access to video memory and a bunch of GPU registers, the GPU read/write access to (a portion of) main memory, and everyone a headache because the latency for all these transactions is even worse than memory latency because the signals have to go out of the chip, into the slot, travel a bit across the mainboard then get to someplace in the CPU about a week later (or that's how it feels compared to the CPU/GPU speeds anyway). The bandwidth is decent though – up to about 8GB/s (theoretical) peak aggregate bandwidth across the 16-lane PCIe 2.0 connections that most GPUs use right now, so between half and a third of the aggregate CPU memory bandwidth; that's a usable ratio. And unlike earlier standards like AGP, this is a symmetrical point-to-point link – that bandwidth goes both directions; AGP had a fast channel from the CPU to the GPU, but not the other way round.

Some final memory bits and pieces

Honestly, we're *very very* close to actually seeing 3D commands now! So close you can almost taste them. But there's one more thing we need to get out of the way first. Because now we have two kinds of memory – (local) video memory and mapped system memory. One is about a day's worth of travel to the north, the other is a week's journey to the south along the PCI Express highway. Which road do we pick?

The easiest solution: Just add an extra address line that tells you which way to go. This is simple, works just fine and has been done plenty of times. Or maybe you're on a unified memory architecture, like some game consoles (but not PCs). In that case, there's no choice; there's just *the memory*, which is where you go, period. If you want something fancier, you add a MMU (memory management unit), which gives you a fully virtualized address space and allows you to pull nice tricks like having frequently accessed parts of a texture in video memory (where they're fast), some other parts in system memory, and most of it not mapped at all – to be conjured up from thin air, or, more usually, by a magic disk read that will only take about 50 years or so – and by the way, this is not hyperbole; if you stay with the "memory access = 1 day" metaphor, that's really how long a single HD read takes. A quite fast one at that. Disks suck. But I digress.

So, MMU. It also allows you to defragment your video memory address space without having to actually copy stuff around when you start running out of video memory. Nice thing, that. And it makes it much easier to have multiple processes share the same GPU. It's definitely allowed to have one, but I'm not actually sure if it's a requirement or not, even though it's certainly really nice to have (anyone care to help me out here? I'll update the article if I get clarification on this, but tbh right now I just can't be arsed to look it up). Anyway, a MMU/virtual memory is not really something you can just add on the side (not in an architecture with caches and memory consistency concerns anyway), but it really isn't specific to any particular stage – I have to mention it somewhere, so I just put it here.

There's also a DMA engine that can copy memory around without having to involve any of our precious 3D hardware/shader cores. Usually, this can at least copy between system memory and video memory (in both directions). It often can also copy from video memory to video memory (and if you have to do any VRAM defragmenting, this is a useful thing to have). It usually can't do system memory to system memory copies, because this is a GPU, not a memory copying unit – do your system memory copies on the CPU where they don't have to pass through PCIe in both directions!

Update: I've drawn a picture (http://www.farbrausch.de/~fg/gpu/gpu_memory.jpg) (link since this layout is too narrow to put big diagrams in the text). This also shows some more details – by now your GPU has multiple memory controllers, each of which controls multiple memory banks, with a fat hub in the front. Whatever it takes to get that bandwidth. :)

Okay, checklist. We have a command buffer prepared on the CPU. We have the PCIe host interface, so the CPU can actually tell us about this, and write its address to some register. We have the logic to turn that address into a load that will actually return data – if it's from system memory it goes through PCIe, if we decide we'd rather have the command buffer in video memory, the KMD can set up a DMA transfer so neither the CPU nor the shader cores on the GPU need to actively worry about it. And then we can get the data from our copy in video memory through the memory subsystem. All paths accounted for, we're set and finally ready to look at some commands!

At long last, the command processor!

Our discussion of the command processor starts, as so many things do these days, with a single word:

"Buffering..."

As mentioned above, both of our memory paths leading up to here are high-bandwidth but also high-latency. For most later bits in the GPU pipeline, the method of choice to work around this is to run lots of independent threads. But in this case, we only have a single command processor that needs to chew through our command buffer in order (since this command buffer contains things such as state changes and rendering commands that need to be executed in the right sequence). So we do the next best thing: Add a large enough buffer and prefetch far enough ahead to avoid hiccups.

From that buffer, it goes to the actual command processing front end, which is basically a state machine that knows how to parse commands (with a hardware-specific format). Some commands deal with 2D rendering operations – unless there's a separate command processor for 2D stuff and the 3D frontend never even sees it. Either way, there's still dedicated 2D hardware hidden on modern GPUs, just as there's a VGA chip somewhere on that die that still supports text mode, 4-bit/pixel bit-plane modes, smooth scrolling and all that stuff. Good luck finding any of that on the die without a microscope. Anyway, that stuff exists, but henceforth I shall not mention it again. :) Then there's commands that actually hand some primitives to the 3D/shader pipe, woo-hoo! I'll take about them in upcoming parts. There's also commands that go to the 3D/shader pipe but never render anything, for various reasons (and in various pipeline configurations); these are up even later.

Then there's commands that change state. As a programmer, you think of them as just changing a variable, and that's basically what happens. But a GPU is a massively parallel computer, and you can't just change a global variable in a parallel system and hope that everything works out OK – if you can't guarantee that everything will work by virtue of some invariant you're enforcing, there's a bug and you will hit it eventually. There's several popular methods, and basically all chips use different methods for different types of state.

Whenever you change a state, you require that all pending work that might refer to that state be finished (i.e. basically a partial pipeline flush). Historically, this is how graphics chips handled most state changes – it's simple and not that costly if you have a low number of batches, few triangles and a short pipeline. Alas, batch and triangle counts have gone up and pipelines have gotten long, so the cost for this type of approach has shot up. It's still alive and kicking for stuff that's either changed infrequently (a dozen partial pipeline flushes aren't that big a deal over the course of a whole frame) or just too expensive/difficult to implement with more specific schemes though.

You can make hardware units completely stateless. Just pass the state change command through up to the stage that cares about it; then have that stage append the current state to everything it sends downstream, every cycle. It's not stored anywhere – but it's always around, so if some pipeline stage wants to look at a few bits in the state it can, because they're passed in (and then passed on to the next stage). If your state happens to be just a few bits, this is fairly cheap and practical. If it happens to be the full set of active textures along with texture sampling state, not so much.

Sometimes storing just one copy of the state and having to flush every time that stage changes serializes things too much, but things would really be fine if you had two copies (or maybe four?) so your state-setting frontend could get a bit ahead. Say you have enough registers ("slots") to store two versions of every state, and some active job references slot 0. You can safely modify slot 1 without stopping that job, or otherwise interfering with it at all. Now you don't need to send the whole state around through the pipeline – only a single bit per command that selects whether to use slot 0 or 1. Of course, if both slot 0 and 1 are busy by the time a state change command is encountered, you still have to wait, but you can get one step ahead. The same technique works with more than two slots.

For some things like sampler or `texture` Shader Resource View state, you could be setting very large numbers of them at the same time, but chances are you aren't. You don't want to reserve state space for 2*128 active textures just because you're keeping track of 2 in-flight state sets so you might need it. For such cases, you can use a kind of register renaming scheme – have a pool of 128 physical texture descriptors. If someone actually needs 128 textures in one shader, then state changes are gonna be slow. (Tough break). But in the more likely case of an app using less than 20 textures, you have quite some headroom to keep multiple versions around.

This is not meant to be a comprehensive list – but the main point is that something that looks as simple as changing a variable in your app (and even in the UMD/KMD and the command buffer for that matter!) might actually need a nontrivial amount of supporting hardware behind it just to prevent it from slowing things down.

Synchronization

Finally, the last family of commands deals with CPU/GPU and GPU/GPU synchronization.

Generally, all of these have the form “if event X happens, do Y”. I’ll deal with the “do Y” part first – there’s two sensible options for what Y can be here: it can be a push-model notification where the GPU yells at the CPU to do something *right now* (“Oi! CPU! I’m entering the vertical blanking interval on display 0 right now, so if you want to flip buffers without tearing, this would be the time to do it!”), or it can be a pull-model thing where the GPU just memorizes that something happened and the CPU can later ask about it (“Say, GPU, what was the most recent command buffer fragment you started processing?” – “Let me check... sequence id 303.”). The former is typically implemented using interrupts and only used for infrequent and high-priority events because interrupts are fairly expensive. All you need for the latter is some CPU-visible GPU registers and a way to write values into them from the command buffer once a certain event happens.

Say you have 16 such registers. Then you could assign `currentCommandBufferSeqId` to register 0. You assign a sequence number to every command buffer you submit to the GPU (this is in the KMD), and then at the start of each command buffer, you add a “If you get to this point in the command buffer, write to register 0”. And voila, now we know which command buffer the GPU is currently chewing on! And we know that the command processor finishes commands strictly in sequence, so if the first command in command buffer 303 was executed, that means all command buffers up to and including sequence id 302 are finished and can now be reclaimed by the KMD, freed, modified, or turned into a cheesy amusement park.

We also now have an example of what X could be: “if you get here” – perhaps the simplest example, but already useful. Other examples are “if all shaders have finished all texture reads coming from batches before this point in the command buffer” (this marks safe points to reclaim texture/render target memory), “if rendering to all active render targets/UAVs has completed” (this marks points at which you can actually safely use them as textures), “if all operations up to this point are fully completed”, and so on.

Such operations are usually called “fences”, by the way. There’s different methods of picking the values you write into the status registers, but as far as I am concerned, the only sane way to do it is to use a sequential counter for this (probably stealing some of the bits for other information). Yeah, I’m really just dropping that one piece of random information without any rationale whatsoever here, because I think you should know. I might elaborate on it in a later blog post (though not in this series) :).

So, we got one half of it – we can now report status back from the GPU to the CPU, which allows us to do sane memory management in our drivers (notably, we can now find out when it’s safe to actually reclaim memory used for vertex buffers, command buffers, textures and other resources). But that’s not all of it – there’s a puzzle piece missing. What if we need to synchronize purely on the GPU side, for example? Let’s go back to the render target example. We can’t use that as a texture until the rendering is actually finished (and some other steps have taken place – more details on that once I get to the texturing units). The solution is a “wait”-style instruction: “Wait until register M contains value N”. This can either be a compare for equality, or less-than (note you need to deal with wraparounds here!), or more fancy stuff – I’m just going with equals for simplicity. This allows us to do the render target sync before we submit a batch. It also allows us to build a full GPU flush operation: “Set register 0 to ++seqId if all pending jobs finished” / “Wait until register 0 contains seqId”. Done and done. GPU/GPU synchronization: solved – and until the introduction of DX11 with Compute Shaders that have another type of more fine-grained synchronization, this was usually the *only* synchronization mechanism you had on the GPU side. For regular rendering, you simply don’t need more.

By the way, if you can write these registers from the CPU side, you can use this the other way too – submit a partial command buffer including a wait for a particular value, and then change the register from the CPU instead of the GPU. This kind of thing can be used to implement D3D11-style multithreaded rendering where you can submit a batch that references vertex/index buffers that are still locked on the CPU side (probably being written to by another thread). You simply stuff the wait just in front of the actual render call, and then the CPU can change the contents of the register once the vertex/index buffers are actually unlocked. If the GPU never got that far in the command buffer, the wait is now a no-op; if it did, it spend some (command processor) time spinning until the data was actually there. Pretty nifty, no? Actually, you can implement this kind of thing even without CPU-writeable status registers if you can modify the command buffer after you submit it, as long as there’s a command buffer “jump” instruction. The details are left to the interested reader :)

Of course, you don’t necessarily need the set register/wait register model; for GPU/GPU synchronization, you can just as simply have a “rendertarget barrier” instruction that makes sure a rendertarget is safe to use, and a “flush everything” command. But I like the set register-style model more because it kills two birds (back-reporting of in-use resources to the CPU, and GPU self-synchronization) with one well-designed stone.

Update: Here, I've drawn a diagram (http://www.farbrausch.de/~fg/gpu/command_processor.jpg) for you. It got a bit convoluted so I'm going to lower the amount of detail in the future. The basic idea is this: The command processor has a FIFO in front, then the command decode logic, execution is handled by various blocks that communicate with the 2D unit, 3D front-end (regular 3D rendering) or shader units directly (compute shaders), then there's a block that deals with sync/wait commands (which has the publicly visible registers I talked about), and one unit that handles command buffer jumps/calls (which changes the current fetch address that goes to the FIFO). And all of the units we dispatch work to need to send us back completion events so we know when e.g. textures aren't being used anymore and their memory can be reclaimed.

Closing remarks

Next step down is the first one doing any actual rendering work. Finally, only 3 parts into my series on GPUs, we actually start looking at some vertex data! (No, no triangles being rasterized yet. That will take some more time).

Actually, at this stage, there's already a fork in the pipeline; if we're running compute shaders, the next step would already be ... running compute shaders. But we aren't, because compute shaders are a topic for later parts! Regular rendering pipeline first.

Small disclaimer: Again, I'm giving you the broad strokes here, going into details where it's necessary (or interesting), but trust me, there's a lot of stuff that I dropped for convenience (and ease of understanding). That said, I don't think I left out anything really important. And of course I might've gotten some things wrong. If you find any bugs, tell me!

Until the next part...

From → Coding, Graphics Pipeline

13 Comments

1. atyuwen permalink

Great article! Is there any book or paper that introduces this kind of stuff detailedly and systematically?

Reply

- [fgiesen permalink](#)

"Real-Time Rendering" (2nd and 3rd Editions) have a chapter each dedicated to graphics hardware that discusses some existing GPU architectures – at a lower level of detail than what I'm doing here, though. For every generation of GPU you'll find some presentations and white papers that explain at least the broad strokes of the architecture – look among Siggraph papers for the last few years, for example. Stay away from marketing blurb and most hardware review sites – most of that is a mish-mash between facts, extrapolation and pure fiction, and it's hard to see what is what. I'm not aware of any book that has an in-depth explanation of GPU architecture, but there is Hennessy and Patterson's "Computer Architecture: A Quantitative Approach" which covers some of this ground (in particular, everything about memory architecture and pipelines is readily applicable, as are the chapters about multiprocessing and thread-level parallelism).

Reply

- [atyuwen permalink](#)

Really appreciate your reply. "Real-time Rendering" is also my favorite book on graphics.

2. skp permalink

> If your state happens to be just a few bits, this isn't fairly cheap and practical.

Should that be 'is' instead of "isn't" ?

Reply

- [fgiesen permalink](#)

It should. Thanks!

Reply

3. ridershen permalink

I follow the diagram link, but haven't found your diagram.

Reply

4. Pieter Kockx permalink

Great article! Small typo that confused me: divide latency (seconds) by clock frequency (cycles/second) should be multiply!

Reply

- [fgiesen permalink](#)

Thanks. Indeed!

Reply

Trackbacks & Pingbacks

1. Geeks3D Programming Links – July 01, 2011 - 3D Tech News, Pixel Hacking, Data Visualization and 3D Programming - Geeks3D.com
2. A trip through the Graphics Pipeline 2011, part 4 « The ryg blog
3. (Updated) 3D Graphics Pipeline Explained - 3D Tech News, Pixel Hacking, Data Visualization and 3D Programming - Geeks3D.com
4. A trip through the Graphics Pipeline 2011: Index « The ryg blog
5. A trip through the Graphics Pipeline 2011, part 9 « The ryg blog

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 3

July 3, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

At this point, we’ve sent draw calls down from our app all the way through various driver layers and the command processor; now, finally we’re actually going to do some graphics processing on it! In this part, I’ll look at the vertex pipeline. But before we start...

Have some Alphabet Soup!

We’re now in the 3D pipeline proper, which in turn consists of several stages, each of which does one particular job. I’m gonna give names to all the stages I’ll talk about – mostly sticking with the “official” D3D10/11 names for consistency – plus the corresponding acronyms. We’ll see all of these eventually on our grand tour, but it’ll take a while (and several more parts) until we see most of them – seriously, I made a small outline of the ground I want to cover, and this series will keep me busy for at least 2 weeks! Anyway, here goes, together with a one-sentence summary of what each stage does.

IA — Input Assembler. Reads index and vertex data.

VS — Vertex shader. Gets input vertex data, writes out processed vertex data for the next stage.

PA — Primitive Assembly. Reads the vertices that make up a primitive and passes them on.

HS — Hull shader; accepts patch primitives, writes transformed (or not) patch control points, inputs for the domain shader, plus some extra data that drives tessellation.

TS — Tessellator stage. Creates vertices and connectivity for tessellated lines or triangles.

DS — Domain shader; takes shaded control points, extra data from HS and tessellated positions from TS and turns them into vertices again.

GS — Geometry shader; inputs primitives, optionally with adjacency information, then outputs different primitives. Also the primary hub for...

SO — Stream-out. Writes GS output (i.e. transformed primitives) to a buffer in memory.

RS — Rasterizer. Rasterizes primitives.

PS — Pixel shader. Gets interpolated vertex data, outputs pixel colors. Can also write to UAVs (unordered access views).

OM — Output merger. Gets shaded pixels from PS, does alpha blending and writes them back to the backbuffer.

CS — Compute shader. In its own pipeline all by itself. Only input is constant buffers+thread ID; can write to buffers and UAVs.

And now that that’s out of the way, here’s a list of the various data paths I’ll be talking about, in order: (I’ll leave out the IA, PA, RS and OM stages in here, since for our purposes they don’t actually do anything *to* the data, they just rearrange/reorder it – i.e. they’re essentially glue)

1. VS→PS: Ye Olde Programmable Pipeline. In D3D9, this was all you got. Still the most important path for regular rendering by far. I’ll go through this from beginning to end then double back to the fancier paths once I’m done.
2. VS→GS→PS: Geometry Shading (new with D3D10).
3. VS→HS→TS→DS→PS, VS→HS→TS→DS→GS→PS: Tessellation (new in D3D11).
4. VS→SO, VS→GS→SO, VS→HS→TS→DS→GS→SO: Stream-out (with and without tessellation).
5. CS: Compute. New in D3D11.

And now that you know what’s coming up, let’s get started on vertex shaders!

Input Assembler stage

The very first thing that happens here is loading indices from the index buffer – if it's an indexed batch. If not, just pretend it was an identity index buffer (0 1 2 3 4 ...) and use that as index instead. If there is an index buffer, its contents are read from memory at this point – not directly though, the IA usually has a data cache to exploit locality of index/vertex buffer access. Also note that index buffer reads (in fact, all resource accesses in D3D10+) are bounds checked; if you reference elements outside the original index buffer (for example, issue a `DrawIndexed` with `IndexCount == 6` from a 5-index buffer) all out-of-bounds reads return zero. Which (in this particular case) is completely useless, but well-defined. Similarly, you can issue a `DrawIndexed` with a NULL index buffer set – this behaves the same way as if you had an index buffer of size zero set, i.e. all reads are out-of-bounds and hence return zero. With D3D10+, you have to work some more to get into the realm of undefined behavior. :)

Once we have the index, we have all we need to read both per-vertex and per-instance data (the current instance ID is just another counter, fairly straightforward, at this stage anyway) from the input vertex streams. This is fairly straightforward – we have a declaration of the data layout; just read it from the cache/memory and unpack it into the float format that our shader cores want for input. However, this read isn't done immediately; the hardware is running a cache of shaded vertices, so that if one vertex is referenced by multiple triangles (and in a fully regular closed triangle mesh, each vertex will be referenced by about 6 tris!) it doesn't need to be shaded every time – we just reference the shaded data that's already there!

Vertex Caching and Shading

Note: The contents of this section are, in part, guesswork. They're based on public comments made by people "in the know" about current GPUs, but that only gives me the "what", not the "why", so there's some extrapolation here. Also, I'm simply guessing some of the details here. That said, I'm not talking completely out of my ass here – I'm confident that what I'm describing here is both reasonable and works (in the general sense), I just can't guarantee that it's actually that way in real HW or that I didn't miss any tricky details. :)

Anyway. For a long time (up to and including the shader model 3.0 generation of GPUs), vertex and pixel shaders were implemented with different units that had different performance trade-offs, and vertex caches were a fairly simple affair: usually just a FIFO for a small number (think one or two dozen) of vertices, with enough space for a worst-case number of output attributes, using the vertex index as a tag. As said, fairly straightforward stuff.

And then unified shaders happened. If you unify two types of shaders that used to be different, the design is necessarily going to be a compromise. So on the one hand, you have vertex shaders, which (at that time) touched maybe up to 1 million vertices a frame in normal use. On the other hand you had pixel shaders, which at 1920×1200 need to touch *at least* 2.3 million pixels a frame *just to fill the whole screen once* – and a lot more if you want to render anything interesting. So guess which of the two units ended up pulling the short straw?

Okay, so here's the deal: instead of the vertex shader units of old that shaded more or less one vertex at a time, you now have a huge beast of a unified shader unit that's designed for maximum throughput, not latency, and hence wants large batches of work (How large? Right now, the magic number seems to be between 16 and 64 vertices shaded in one batch).

So you need between 16-64 vertex cache misses until you can dispatch one vertex shading load, if you don't want to shade inefficiently. But the whole FIFO thing doesn't really play ball with this idea of batching up vertex cache misses and shading them in one go. The problem is this: if you shade a whole batch of vertices at once, that means you can only actually start assembling triangles once all those vertices have finished shading. At which point you've just added a whole batch (let's just say 32 here and in the following) of vertices to the end of the FIFO, which means 32 old vertices now fell out – but each of these 32 vertices might've been a vertex cache hit for one of the triangles in the current batch we're trying to assemble! Uh oh, that doesn't work. Clearly, we can't actually count the 32 oldest verts in the FIFO as vertex cache hits, because by the time we want to reference them they'll be gone! Also, how big do we want to make this FIFO? If we're shading 32 verts in a batch, it needs to be at least 32 entries large, but since we can't use the 32 oldest entries (since we'll be shifting them out), that means we'll effectively start with an empty FIFO on every batch. So, make it bigger, say 64 entries? That's pretty big. And note that every vertex cache lookup involves comparing the tag (vertex index) against all tags in the FIFO – this is fully parallel, but it also a power hog; we're effectively implementing a fully associative cache here. Also, what do we do between dispatching a shading load of 32 vertices and receiving results – just wait? This shading will take a few hundred cycles, waiting seems like a stupid idea! Maybe have two shading loads in flight, in parallel? But now our FIFO needs to be at least 64 entries long, and we can't count the last 64 entries as vertex cache hits, since they'll be

shifted out by the time we receive results. Also, one FIFO vs. lots of shader cores? **Amdahl's law**

(http://en.wikipedia.org/wiki/Amdahl%27s_law) still holds – putting one strictly serial component in a pipeline that's otherwise completely parallel is a surefire way to make it the bottleneck.

This whole FIFO thing really doesn't adapt well to this environment, so, well, just throw it out. Back to the drawing board. What do we actually want to do? Get a decently-sized batch of vertices to shade, and not shade vertices (much) more often than necessary.

So, well, keep it simple: Reserve enough buffer space for 32 vertices (=1 batch), and similarly cache tag space for 32 entries. Start with an empty "cache", i.e. all entries invalid. For every primitive in the index buffer, do a lookup on all the indices; if it's a hit in the cache, fine. If it's a miss, allocate a slot in the current batch and add the new index to the cache tag array. Once we don't have enough space left to add a new primitive anymore, dispatch the whole batch for vertex shading, save the cache tag array (i.e. the 32 indices of the vertices we just shaded), and start setting up the next batch, again from an empty cache – ensuring that the batches are completely independent.

Each batch will keep a shader unit busy for some while (probably at least a few hundred cycles!). But that's no problem, because we got plenty of them – just pick a different unit to execute each batch! Presto parallelism. We'll eventually get the results back. At which point we can use the saved cache tags and the original index buffer data to assemble primitives to be sent down the pipeline (this is what "primitive assembly" does, which I'll cover in the later part).

By the way, when I say "get the results back", what does that mean? Where do they end up? There's two major choices: 1. specialized buffers or 2. some general cache/scratchpad memory. It used to be 1), with a fixed organization designed around vertex data (with space for 16 float4 vectors of attributes per vertex and so on), but lately GPUs seem to be moving towards 2), i.e. "just memory". It's more flexible, and has the distinct advantage that you can use this memory for other shader stages, whereas things like specialized vertex caches are fairly useless for the pixel shading or compute pipeline, to give just one example.

Update: And here's a **picture** (http://www.farbrausch.de/~fg/gpu/vertex_shade.jpg) of the vertex shading dataflow as described so far.

Shader Unit internals

Short versions: It's pretty much what you'd expect from looking at disassembled HLSL compiler output (`fxc /dumpbin` is your friend!). Guess what, it's just processors that are *really good* at running that kind of code, and the way that kind of thing is done in hardware is building something that eats something fairly close to shader bytecode, in spirit anyway. And unlike the stuff that I've been talking about so far, it's fairly well documented too – if you're interested, just check out conference presentations from AMD and NVidia or read the documentation for the CUDA/Stream SDKs.

Anyway, here's the executive summary: fast ALU mostly built around a FMAC (Floating Multiply-ACcumulate) unit, some HW support for (at least) reciprocal, reciprocal square root, log2, exp2, sin, cos, optimized for high throughput and high density not low latency, running a high number of threads to cover said latency, fairly small number of registers per thread (since you're running so many of them!), very good at executing straight-line code, bad at branches (especially if they're not coherent).

All that is common to pretty much all implementations. There's some differences, too; AMD hardware used to stick directly with the 4-wide SIMD implied by the HLSL/GLSL and shader bytecode (even though they seem to be moving away from that lately), while NVidia decided to rather turn the 4-way SIMD into scalar instructions a while back. Again though, all that's on the Web already!

What's interesting to note though is the *differences* between the various shader stages. The short version is that really are rather few of them; for example, all the arithmetic and logic instructions are exactly the same across all stages. Some constructs (like derivative instructions and interpolated attributes in pixel shaders) only exist in some stages; but mostly, the differences are just what kind (and format) of data are passed in and out.

There's one special bit related to shaders though that's a big enough subject to deserve a part on its own. That bit is texture sampling (and texture units). Which, it turns out, will be our topic next time! See you then.

Closing remarks

Again, I repeat my disclaimer from the “Vertex Caching and Shading” section: Part of that is conjecture on my part, so take it with a grain of salt. Or maybe a pound. I don’t know.

I’m also not going into any detail on how scratch/cache memory is managed; the buffer sizes depend (primarily) on the size of batches you process and the number of vertex output attributes you expect. Buffer sizing and management is *really* important for performance, but I can’t meaningfully explain it here, nor do I want to; while interesting, this stuff is very specific to whatever hardware you’re talking about, and not really very insightful.

From → Coding, Graphics Pipeline

13 Comments

1. Won Chun permalink

Keep it up! These articles are great. I used to know this nitty gritty way better years ago, but things have changed so much since then and this has been a great way to catch up.

So I thought one of the reasons why the vertex caches changed was not just because of the massive parallelism of the shaders, but also because primitive assembly became a bottleneck, and 1 tri/clock was no longer sufficient. I could be wrong or redundant here, of course.

Transparent post-transform vertex caching is pretty cool, but I always thought that it would have been a big win to dispose of the associative lookups, and use an explicit style (kind of like an extension to the old generalized strips approach). Without associative lookups, you could have much larger, explicitly indexed caches.

The way it would work, is that instead of having index buffers, you’d have a list of back-references to previously seen vertices. These references could be small since they index the cache, rather than the vertex buffer; they could be 8-bit, instead of 16 or 32-bit. You would need a few escape codes, like to index a never-seen vertex, a previously-seen vertex out of the cache, and maybe a primitive restart code.

Maybe this makes more sense for mobile (many implementations don’t really have post-transform caches and rely only on strips, maybe because of the power issue) than desktop.

Reply

o fgiesen permalink

“So I thought one of the reasons why the vertex caches changed was not just because of the massive parallelism of the shaders, but also because primitive assembly became a bottleneck, and 1 tri/clock was no longer sufficient. I could be wrong or redundant here, of course.”

Yep, that’s what I was hinting with the “the FIFO is a serial part of the pipeline” bit; basically PA needs to have the vertex data for all of the primitive somewhere, and if that’s place is your FIFO, then you can have only one PA (because there’s just one FIFO)! You could have multiple FIFOs in theory, but then you’d end up re-shading vertices for every FIFO individually, which is totally beside the point for a vertex cache. So yeah, just block it and dispatch/work on a chunk basis downstream. Way easier (and more scalable!) overall.

“Transparent post-transform vertex caching is pretty cool, but I always thought that it would have been a big win to dispose of the associative lookups, and use an explicit style”

So basically a two-level indexing scheme; I’m not sure. It’s certainly a more efficient way to drive the pipeline in general, but interface compatibility is a big deal. No matter what they do, they’ll still need to support regular indexed primitives efficiently, since that’s what everyone uses up until now. At this point there’s just significant friction involved in changing this; we might get there eventually if it becomes a significant issue, but right now that doesn’t seem to be the case.

Also note that e.g. with a 16-wide shader unit, you want to shoot for having batches at least as large as 16, but you can obviously shoot for wider multiples too: 32, 48, 64 and so on. This gives you better hit rates (this whole setup thing I described restarts from scratch in every block!) but also needs more memory for buffering and has higher granularity. If you go for a fully associative cache, there’ll be a fixed limit of how large your maximum vertex shade batch can be (which is the size of that cache). But the better way to make it scale really wide is to just not be fully associative, and make the cache set-associative instead; now every vertex index has a couple places it can go in the cache. If you encounter an index that would go into a set that’s already full, you can’t add it to the current block anymore, so you have to dispatch what might be a partial block. Note that for triangles your cache needs to be at least 3-way associative so you can make progress – all vertex indices in a tri might map to the same set!

I don't think anyone currently does this though; for one, very wide batches need lots of buffering for output, and buffer space is limited, so larger batches mean you can run fewer of them in parallel. If you're running too few in parallel, you can't cover the latency and end up waiting. More importantly, I'm talking about D3D11 here; the worst-case primitive that needs to be supported is a 32-control point patch, which absolutely positively needs at least a 32-vertex batch in order to assemble a single primitive! (You need to be able to assemble at least one primitive at all times, or you could get stuck indefinitely). Which would also mean you need at least a 32-way associative cache in the set associative model!

"Maybe this makes more sense for mobile (many implementations don't really have post-transform caches and rely only on strips, maybe because of the power issue) than desktop."

The thing with strips is that you basically get them for free. In PA, you need buffers for at least 3 vertices (a single triangle) no matter what you do. But once you buffer the last 3 vertices, you get triangle strips effectively for free. A proper post-transform vertex cache, on the other hand, takes up area, power and design effort. If you don't care much about triangle throughput, it's not worth it.

Reply

- o **zeuxcg** [permalink](#)

As far as I know, the scheme more or less like yours was used in old GeForces ("old" means before GFFX, I think – maybe even before GF3). I also vaguely remember that dynamic index buffers were slow, which is related, of course.

Still, the current way is better – it opens ways to do things differently (i.e. it's hard to do backreferences in a parallel way, the same way as it's hard to do traditional FIFO cache), it's conceptually simple and I don't think it wastes any considerable amount of energy/die space.

Reply

- o **fgiesen** [permalink](#)

If this was the case, they can't have used it for long; not sure about GF1/2, but GF3 had a FIFO scheme like what I described, as did the Radeons at the time.

GF3 definitely had slow *static* index buffers (dynamic was fine). IIRC that was because they didn't support index buffers in hardware at all – index data had to come from the command buffer, so the driver needed to `memcpy` index data around.

2. Corbin Simpson [permalink](#)

A great article.

For future reference, you can point people at The X.Org Foundation's repository of documentation, which is all publicly available at <http://www.x.org/docs/>. For example, AMD's r600/r700 shader documentation, at the assembly level, is publicly viewable at <http://www.x.org/docs/AMD/r600isa.pdf>. Take it easy!

Reply

3. Ignacio [permalink](#)

I can confirm that your guesses are correct, at least for NVIDIA hardware. Mark Kilgard wrote a bit about the block-based vertex cache in modern GPUs here:

http://www.slideshare.net/Mark_Kilgard/using-vertex-bufferobjectsowell

It'd be interesting to do some numbers and compare the efficiency of FIFO vs. block caches. For the same size you would find that FIFOs always do a lot better. This is even more so when using primitives with more than 3 vertices. For example, the efficiency when rendering 16×16 patches is terrible, and since vertices are shared by many more patches than triangles...

Reply

4. TomF [permalink](#)

Just to note that every architecture now turns DXasm code into scalar code. "add r0.xyz, r1, r2" turns into three separate instructions, and r0.x, r0.y and r0.z are three completely separate registers. Even the older AMD/ATI architectures that still had VLIW5 or VLIW4 started by scalarising the shader, doing clever things with it, and then re-VLIW'd in cunning ways completely unrelated to how DXasm arranged things. (see the very excellent talk by Norm Rubin "Issues and challenges in compiling for graphics processors", <http://dl.acm.org/citation.cfm?id=1356088>)

What does this mean in practical terms? Well first – use the correct write masks – don't rely on the dead-code-elimination of the compiler. Second – nothing is free any more – if you don't need a value, don't calculate it, even if it's "just" the fourth channel. Finally, if you're counting instructions, you always need to multiply by the number of destination targets. "mul r1, r2, r3" is four times as expensive as "mul r1.x, r2, r3". (obvious exceptions for things like dp4 and so on)

Reply

5. Bala permalink

Hi,

You postings are really interesting. I am using Amd graphics processors have rv730 architecture. 4690. Do you have a good reference docs or material which explains the interernals of this architecture.

Reply

Trackbacks & Pingbacks

1. (Updated) 3D Graphics Pipeline Explained - 3D Tech News, Pixel Hacking, Data Visualization and 3D Programming - Geeks3D.com
2. A trip through the Graphics Pipeline 2011: Index « The ryg blog
3. A trip through the Graphics Pipeline 2011, part 10 « The ryg blog
4. A trip through the Graphics Pipeline 2011, part 12 « The ryg blog
5. A trip through the Graphics Pipeline 2011, part 13 « The ryg blog

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 4

July 4, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Welcome back. Last part was about vertex shaders, with some coverage of GPU shader units in general. Mostly, they’re just vector processors, but they have access to one resource that doesn’t exist in other vector architectures: Texture samplers. They’re an integral part of the GPU pipeline and are complicated (and interesting!) enough to warrant their own article, so here goes.

Texture state

Before we start with the actual texturing operations, let’s have a look at the API state that drives texturing. In the D3D11 part, this is composed of 3 distinct parts:

1. The sampler state. Filter mode, addressing mode, max anisotropy, stuff like that. This controls how texture sampling is done in a general way.
2. The underlying texture resource. This boils down to a pointer to the raw texture bits in memory. The resource also determines whether it’s a single texture or a texture array, what multisample format the texture has (if any), and the physical layout of the texture bits – i.e. at the resource level, it’s not yet decided how the values in memory are to be interpreted exactly, but their memory layout is nailed down.
3. The shader resource view (SRV for short). This determines how the texture bits are to be interpreted by the sampler. In D3D10+, the resource view links to the underlying resource, so you never specify the resource explicitly.

Most of the time, you will create a texture resource with a given format, let’s say RGBA, 8 bits per component, and then just create a matching SRV. But you can also create a texture as “8 bits per component, typeless” and then have several different SRVs for the same resource that read the underlying data in different formats, e.g. once as UNORM8_SRGB (unsigned 8-bit value in sRGB space that gets mapped to float 0..1) and once as UINT8 (unsigned 8-bit integer).

Creating the extra SRV seems like an annoying extra step at first, but the point is that this allows the API runtime to do all type checking at SRV creation time; if you get a valid SRV back, that means the SRV and resource formats are compatible, and no further type checking needs to be done while that SRV exists. In other words, it’s all about API efficiency here.

Anyway, at the hardware level, what this boils down to is just a bag of state associated with a texture sampling operation – sampler state, texture/format to use, etc. – that needs to get kept somewhere (see [part 2](#) (<https://fgiesen.wordpress.com/2011/07/02/a-trip-through-the-graphics-pipeline-2011-part-2/>) for an explanation of various ways to manage state in a pipelined architecture). So again, there’s various methods, from “pipeline flush every time any state changes” to “just go completely stateless in the sampler and send the full set along with every texture request”, with various options inbetween. It’s nothing you need to worry about – this is the kind of thing where HW architects whip up a cost-benefit analysis, simulate a few workloads and then take whichever method comes out ahead – but it’s worth repeating: as PC programmer, don’t assume the HW adheres to any particular model.

Don’t assume that texture switches are expensive – they might be fully pipelined with stateless texture samplers so they’re basically free. But don’t assume they’re completely free either – maybe they are not fully pipelined or there’s a cap on the maximum number of different sets of texture states in the pipeline at any given time. Unless you’re on a console with fixed hardware (or you hand-optimize your engine for every generation of graphics HW you’re targeting), there’s just no way to tell. So when optimizing, do the obvious stuff – sort by material where possible to avoid unnecessary state changes and the like – which certainly saves you some API work at the very least, and then leave it at that. Don’t do anything fancy based on any particular model of what the HW is doing, because it can (and will!) change in the blink of an eye between HW generations.

Anatomy of a texture request

So, how much information do we need to send along with a texture sample request? It depends on the texture type and which kind of sampling instruction we're using. For now, let's assume a 2D texture. What information do we need to send if we want to do a 2D texture sample with, say, up to 4x anisotropic sampling?

The 2D texture coordinates – 2 floats, and sticking with the D3D terminology in this series, I'm going to call them u/v and not s/t.

The partial derivatives of u and v along the screen "x" direction: $\frac{\partial u}{\partial x}, \frac{\partial v}{\partial x}$.

Similarly, we need the partial derivative in the "y" direction too: $\frac{\partial u}{\partial y}, \frac{\partial v}{\partial y}$.

So, that's 6 floats for a fairly pedestrian 2D sampling request (of the `SampleGrad` variety) – probably more than you thought. The 4 gradient values are used both for mipmap selection and to choose the size and shape of the anisotropic filtering kernel. You can also use texture sampling instructions that explicitly specify a mipmap level (in HLSL, that would be `SampleLevel`) – these don't need the gradients, just a single value containing the LOD parameter, but they also can't do anisotropic filtering – the best you'll get is trilinear! Anyway, let's stay with those 6 floats for a while. That sure seems like a lot. Do we really need to send them along with every texture request?

The answer is: depends. In everything but Pixel Shaders, the answer is yes, we really have to (if we want anisotropic filtering that is). In Pixel Shaders, turns out we don't; there's a trick that allows Pixel Shaders to give you gradient instructions (where you can compute some value and then ask the hardware "what is the approximate screen-space gradient of this value?"), and that same trick can be employed by the texture sampler to get all the required partial derivatives just from the coordinates. So for a PS 2D "sample" instruction, you really only need to send the 2 coordinates which imply the rest, provided you're willing to do some more math in the sampler units.

Just for kicks: What's the worst-case number of parameters required for a single texture sample? In the current D3D11 pipeline, it's a `SampleGrad` on a Cubemap array. Let's see the tally:

3D texture coordinates – u, v, w: 3 floats.

Cubemap array index: one int (let's just bill that at the same cost as a float here).

Gradient of (u,v,w) in the screen x and y directions: 6 floats.

For a total of 10 values *per pixel sampled* – that's 40 bytes if you actually store it like that. Now, you might decide that you don't need full 32 bits for all of this (it's probably overkill for the array index and gradients), but it's still a lot of data to be sending around.

In fact, let's check what kind of bandwidth we're talking about here. Let's assume that most of our textures are 2D (with a few cubemaps thrown in), that most of our texture sampling requests come from the Pixel Shader with little to no texture samples in the Vertex Shader, and that the regular `Sample`-type requests are the most frequent, followed by `SampleLevel` (all of this is pretty typical for actual rendering you see in games). That means the average number of 32-bit floats values sent per pixel will be somewhere between 2 (u+v) and 3 (u+v+w / u+v+lod), let's say 2.5, or 10 bytes.

Assume a medium resolution – say, 1280×720, which is about 0.92 million pixels. How many texture samples does your average game pixel shader have? I'd say at least 3. Let's say we have a modest amount of overdraw, so during the 3D rendering phase, we touch each pixel on the screen roughly twice. And then we finish it off with a few texture-heavy full-screen passes to do post-processing. That probably adds at least another 6 samples per pixel, taking into account that some of that postprocessing will be done at a lower resolution. Add it all up and we have $0.92 * (3*2 + 6)$ = about 11 million texture samples per frame, which at 30 fps is about 330 million a second. At 10 bytes per request, that's 3.3 GB/s *just for texture request payloads*. Lower bound, since there's some extra overhead involved (we'll get to that in a second). Note that I'm *cough* erring "a bit" on the low side with all of these numbers :). An actual modern game on a good DX11 card will run in significantly higher resolution, with more complex shaders than I listed, comparable amount of overdraw or even somewhat less (deferred shading/lighting to the rescue!), higher frame rate, and way more complex postprocessing – go ahead, do a quick back-of-the-envelope calculation how much texture request bandwidth a decent-quality SSAO pass in quarter-resolution with bilateral upsampling takes...

Point being, this whole texture bandwidth thing is not something you can just hand-wave away. The texture samplers aren't part of the shader cores, they're separate units some distance away on the chip, and shuffling multiple gigabytes per second around isn't something that just happens by itself. This is an actual architectural issue – and it's a good thing we don't use `SampleGrad` on Cubemap arrays for everything :)

But who asks for a single texture sample?

The answer is of course: *No one*. Our texture requests are coming from shader units, which we know process somewhere between 16 and 64 pixels / vertices / control points / ... at once. So our shaders won't be sending individual texture samples, they'll dispatch a bunch of them at once. This time, I'll use 16 as the number – simply because the 32 I chose last time is non-square, which just seems weird when talking about 2D texture requests. So, 16 texture requests at once – build that texture request payload, add some command fields at the start so the sampler knows what to do, add some more fields so the sampler knows which texture and sampler state to use (again, see the remarks above on state), and send that off to a texture sampler somewhere.

This will take a while.

No, seriously. Texture samplers have a seriously long pipeline (we'll soon see why); a texture sampling operation takes *way* too long for a shader unit to just sit idle for all that time. Again, say it with me: *throughput*. So what happens is that on a texture sample, a shader unit will just quietly switch to another thread/batch and do some other work, then switch back a while later when the results are there. Works just fine as long as there's enough independent work for the shader units to do!

And once the texture coordinates arrive...

Well, there's a bunch of computations to be done first: (In here and the following, I'm assuming a simple bilinear sample; trilinear and anisotropic take some more work, see below).

If this is a `Sample` or `SampleBias`-type request, calculate texture coordinate gradients first.

If no explicit mip level was given, calculate the mip level to be sampled from the gradients and add the LOD bias if specified. For each resulting sample position, apply the address modes (wrap / clamp / mirror etc.) to get the right position in the texture to sample from, in normalized [0,1] coordinates.

If this is a cubemap, we also need to determine which cube face to sample from (based on the absolute values and signs of the u/v/w coordinates), and do a division to project the coordinates onto the unit cube so they are in the [-1,1] interval. We also need to drop one of the 3 coordinates (based on the cube face) and scale/bias the other 2 so they're in the same [0,1] normalized coordinate space we have for regular texture samples.

Next, take the [0,1] normalized coordinates and convert them into fixed-point pixel coordinates to sample from – we need some fractional bits for the bilinear interpolation.

Finally, from the integer x/y/z and texture array index, we can now compute the address to read texels from. Hey, at this point, what's a few more multiplies and adds among friends?

If you think it sounds bad summed up like that, let me take remind you that this is a simplified view. The above summary doesn't even cover fun issues such as texture borders or sampling cubemap edges/corners. Trust me, it may sound bad now, but if you were to actually write out the code for everything that needs to happen here, you'd be positively horrified. Good thing we have dedicated hardware to do it for us. :) Anyway, we now have a memory address to get data from. And wherever there's memory addresses, there's a cache or two nearby.

Texture cache

Everyone seems to be using a two-level texture cache these days. The second-level cache is a completely bog-standard cache that happens to cache memory containing texture data. The first-level cache is not quite as standard, because it's got additional smarts. It's also smaller than you probably expect – on the order of 4-8kb per sampler. Let's cover the size first, because it tends to come as a surprise to most people.

The thing is this: Most texture sampling is done in Pixel Shaders with mip-mapping enabled, and the mip level for sampling is specifically chosen to make the screen pixel:texel ratio roughly 1:1 – that's the whole point. But this means that, unless you happen to hit the exact same location in a texture again and again, each texture sampling operation will miss about 1 texel on average – the actual measured value with bilinear filtering is around 1.25 misses/request (if you track pixels individually). This value stays more

or less unchanged for a long time even as you change texture cache size, and then drops dramatically as soon as your texture cache is large enough to contain the whole texture (which usually is between a few hundred kilobytes and several megabytes, totally unrealistic sizes for a L1 cache).

Point being, *any* texture cache whatsoever is a massive win (since it drops you down from about 4 memory accesses per bilinear sample down to 1.25). But unlike with a CPU or shared memory for shader cores, there's very little gain in going from say 4k of cache to 16k; we're streaming larger texture data through the cache no matter what.

Second point: Because of the 1.25 misses/sample average, texture sampler pipelines need to be long enough to sustain a full read from memory per sample without stalling. Let me phrase that differently: texture sampler pipes are long enough to not stall for a memory read *even though it takes 400-800 cycles*. That's one seriously long pipeline right there – and it really is a pipeline in the literal sense, handing data from one pipeline register to the next for a few hundred cycles without any processing until the memory read is completed.

So, small L1 cache, long pipeline. What about the “additional smarts”? Well, there’s compressed texture formats. The ones you see on PC – S3TC aka DXTC aka BC1-3, then BC4 and 5 which were introduced with D3D10 and are just variations on DXT, and finally BC6H and 7 which were introduced with D3D11 – are all block-based methods that encode blocks of 4×4 pixels individually. If you decode them during texture sampling, that means you need to be able to decode up to 4 such blocks (if your 4 bilinear sample points happen to land in the worst-case configuration of straddling 4 blocks) per cycle and get a single pixel from each. That, frankly, just sucks. So instead, the 4×4 blocks are decoded when it’s brought into the L1 cache: in the case of BC3 (aka DXT5), you fetch one 128-bit block from texture L2, and then decode that into 16 pixels in the texture cache. And suddenly, instead of having to partially decode up to 4 blocks per sample, you now only need to decode $1.25/(4^2) = \text{about } 0.08$ blocks per sample, at least if your texture access patterns are coherent enough to hit the other 15 pixels you decoded alongside the one you actually asked for :). Even if you only end up using part of it before it goes out of L1 again, that’s still a massive improvement. Nor is this technique limited to DXT blocks; you can handle most of the differences between the >50 different texture formats required by D3D11 in your cache fill path, which is hit about a third as often as the actual pixel read path – nice. For example, things like UNORM sRGB textures can be handled by converting the sRGB pixels into a 16-bit integer/channel (or 16-bit float/channel, or even 32-bit float if you want) in the texture cache. Filtering then operates on that, properly, in linear space. Mind that this does end up increasing the footprint of texels in the L1 cache, so you might want to increase L1 texture size; not because you need to cache more texels, but because the texels you cache are fatter. As usual, it’s a trade-off.

Filtering

And at this point, the actual bilinear filtering process is fairly straightforward. Grab 4 samples from the texture cache, use the fractional positions to blend between them. That’s a few more of our usual standby, the multiply-accumulate unit. (Actually a lot more – we’re doing this for 4 channels at the same time...)

Trilinear filtering? Two bilinear samples and another linear interpolation. Just add some more multiply-accumulates to the pile.

Anisotropic filtering? Now that actually takes some extra work earlier in the pipe, roughly at the point where we originally computed the mip-level to sample from. What we do is look at the gradients to determine not just the area but also the shape of a screen pixel in texel space; if it’s roughly as wide as it is high, we just do a regular bilinear/trilinear sample, but if it’s elongated in one direction, we do several samples across that line and blend the results together. This generates several sample positions, so we end up looping through the full bilinear/trilinear pipeline several times, and the actual way the samples are placed and their relative weights are computed is a closely guarded secret for each hardware vendor; they’ve been hacking at this problem for years, and by now both converged on something pretty damn good at reasonable hardware cost. I’m not gonna speculate what it is they’re doing; truth be told, as a graphics programmer, you just don’t need to care about the underlying anisotropic filtering algorithm as long as it’s not broken and produces either terrible artifacts or terrible slowdowns.

Anyway, aside from the setup and the sequencing logic to loop over the required samples, this does not add a significant amount of computation to the pipe. At this point we have enough multiply-accumulate units to compute the weighted sum involved in anisotropic filtering without a lot of extra hardware in the actual filtering stage. :)

Texture returns

And now we're almost at the end of the texture sampler pipe. What's the result of all this? Up to 4 values (r, g, b, a) per texture sample requested. Unlike texture requests where there's significant variation in the size of requests, here the most common case by far is just the shader consuming all 4 values. Mind you, sending 4 floats back is nothing to sneeze at from a bandwidth point of view, and again you might want to shave bits in some case. If your shader is sampling a 32-bit float/channel texture, you'd better return 32-bit floats, but if it's reading a 8-bit UNORM SRGB texture, 32 bit returns are just overkill, and you can save bandwidth by using a smaller format on the return path.

And that's it – the shader unit now has its texture sampling results back and can resume working on the batch you submitted – which concludes this part. See you again in the next installment, when I talk about the work that needs to be done before we can actually start rasterizing primitives. **Update:** And here's a [picture](http://www.farbrausch.de/~fg/gpu/texture_sample.jpg) (http://www.farbrausch.de/~fg/gpu/texture_sample.jpg) of the texture sampling pipeline, including an amusing mistake that I've fixed in post like a pro!

The usual post-script

This time, no big disclaimers. The numbers I mentioned in the bandwidth example are honestly just made up on the spot since I couldn't be arsed to look up some actual figures for current games :), but other than that, what I describe here should be pretty close to what's on your GPU right now, even though I hand-waved past some of the corner cases in filtering and such (mainly because the details are more nauseating than they are enlightening).

As for texture L1 cache containing decompressed texture data, to the best of my knowledge this is accurate for current hardware. Some older HW would keep some formats compressed even in L1 texture cache, but because of the "1.25 misses/sample for a large range of cache sizes" rule, that's not a big win and probably not worth the complexity. I think that stuff's all gone now.

An interesting bit are embedded/power-optimized graphics chips, e.g. PowerVR; I'll not go into these kinds of chips much in this series since my focus here is on the high-performance parts you see in PCs, but I have some notes about them in the comments for previous parts if you're interested. Anyway, the PVR chips have their own texture compression format that's not block-based and very tightly integrated with their filtering hardware, so I would assume that they do keep their textures compressed even in L1 texture cache (actually, I don't know if they even have a second cache level!). It's an interesting method and probably at a fairly sweet spot in terms of useful work done per area and energy consumed. But I think the "depack to L1 cache" method gives higher throughput overall, and as I can't mention often enough, it's all about throughput on high-end PC GPUs :)

From → Coding, Graphics Pipeline

19 Comments

1. Jocelyn Houle [permalink](#)

Your whole series is very interesting, and filled with much more technical insight that one might expect.

"The above summary doesn't even cover fun issues such as texture borders or sampling cubemap edges/corners. Trust me, it may sound bad now, but if you were to actually write out the code for everything that needs to happen here, you'd be positively horrified."

Amen to that...

I've seen your demoscene experience, and you are now at RAD; you don't seem to have worked at an IHV. Yet, I'm surprised by the depth of your knowledge on hardware, here. Who hinted you at nitty gritty details? Maybe Tom Forsyth, who worked on Larrabee. Anyone else?...

Reply

o fgiesen [permalink](#)

I'll give you the highlights: I've been doing low-level systems/graphics programming for a while, with some compression and compiler work on the side. Some of my friends are GPU HW/SW architects in the embedded space (car navigational systems and such) so I've had a direct line to GPU architects well before I ever met Tom. :) In 2008 I got my degree and went to work in the games industry, doing all kinds of fun work, a lot of it being rendering/optimization work on PS3, Xbox 360 and Wii (sometimes all three at them at once; I needed an extra desk just to stack the Devkits...). Then Jeff hired me to work at RAD on Larrabee, where I arrived in spring 2010 just in time to get a thorough introduction and then a front-row seat for the fireworks (namely it being moved to HPC when all of our contacts at Intel where in the graphics group). So I ended up

helping Sean get Iggy out the door, finishing the half-completed Xbox 360 port and writing new ports for PS3 and Wii, plus a D3D10 renderer on PC – with a lot of optimization work on the side. Meanwhile Intel decided they'd still like us to do some research work on graphics, which is what I've been doing for most of this year so far :)

Reply

2. Jocelyn Houle permalink

Impressive...

I wish you could tell what research work you guys are doing. :-)

Also, about your friends doing GPU HW/SW architecture for car navigators: are they doing shaders on those, these days? Somehow, I doubt they kept up with DX11-class pipeline AMD/NVIDIA have (Matrox couldn't keep up), but sometimes, small teams can do impressive stuff...

Reply

- [fgiesen permalink](#)

Even the GL ES world is all on shaders these days. Hardware that targets 2D/OpenVG instead is more likely to still have a fixed function pixel pipe though (register combiner-style). Vertex shaders are a done deal though; a usable programmable vertex shader unit isn't more complex (in either area, power draw or design/validation effort) than a compliant implementation of the OpenGL fixed-function vertex pipeline: there's position/normal transformation, vertex lighting (for multiple lights!), materials, texture coordinate generation, color/texture matrix, all that crap. If you replace that with a vertex shader unit, you need all the math blocks once instead of multiple times (which is a net win); some buffers get replaced with a register file (more or less a wash), the sequencing/glue logic gets replaced with an instruction decoder/sequencer (a bit more complex so net loss) and validation effort goes from "we have to verify that this HW implements the fixed-function vertex pipeline correctly in all cases" down to "we have to verify that all the math blocks, the IO and the instruction decoder/sequencer work" (absolutely massive win!).

Actually the original HW implementations of OpenGL just used programmable SIMD arrays for both vertex and pixel processing. I'm not even sure whether early PC hardware T&L boards like the GeForce 256 actually implemented the vertex pipe fully in dedicated logic, or whether it was a custom vector processor/DSP that was microcoded. The latter seems more likely, and at that point going to a programmable pipe is a disruptive but incremental change.

Reply

- [Jocelyn Houle permalink](#)

Sure, described this way, VS is an easy switch. But how about efficient PS? I mean, when you open up IEEE-compliant floats, arbitrary shader length (as opposed to 4 texture combiners fixed-function stages), all of the specialized texture fetches, and other things, you end up opening quite a large verification door, I think. And power utilization is king in embedded land. So, if the fixed-function pipeline does the job, is it worth switching architecture?

Are you saying my TomTom could have better 3D capability than my iPod Touch 2nd generation? (OpenGL ES 1.1, not 2.0)

3. fgiesen permalink

I was just talking about VS. I did mention that pixel processing for low-power devices (especially when they don't need to run 3D games) is still doing register combiner-style stuff, and likely to stay that way. Going from a few 8-bit precision blend units to programmable shading is a big deal, and a losing move if your target apps don't actually need it.

Reply

4. shusen permalink

Hi Thanks for the article, I've been looking for something like this for a while.

Also, I wondering if you know how texture was store in video memory (especially 3D texture). I assume it would be some kind of block base scheme (Ziyad S. Hakura's texture cache paper). I hope to get similar performance by implement my own tri-linear interpolation using linear memory (in cuda or openCL) instead of rely on 3D texture.

Reply

- [fgiesen permalink](#)

There's no single answer to that question. Each GPU has its own preferred texture storage formats, which are usually different for 2D and 3D textures (and the array variants), and most GPUs also support multiple layouts for all of them. I didn't go into detail on texture storage because I've written a blog post exclusively on this subject a few months ago (here) which talks about a general family of storage formats that is quite common. Most texture layouts I've seen in practice can be directly expressed within that scheme, and the ones that don't can be closely approximated :).

Most current PC GPUs support linear formats, which are trivial. The various tiled formats vary a lot though. As for volume textures, I believe Nvidia hardware used a straight Morton Order (aka Bit-Interleaving) layout for 3D textures at some point; I'm not sure if that's still the case though, and I also have no clue what AMD uses for 3D textures.

Reply

- **Shusen Liu permalink**

Thanks a lot! That's very helpful. I tried some simple tiling method before I read this, got about 40% of the 3D texture performance. There should be some space for optimization.

5. Barbie permalink

I'm really liking this series. It's got the perfect level of info for somebody who wants to get seriously deep.

I do have one nitpick on this part, though:

"the answer is yes, we really have to (if we want anisotropic filtering that is)"

How do you think you can send down the u,v derivatives on non-PS shaders ? the whole x/y screen position you're supposed to derive on is specific to PS.

Reply

- **fgiesen permalink**

All shader types can use `SampleGrad`. "Sending down" the values is no problem at all – hey, it's just a bunch of floats!

I guess what you mean is how you're supposed to derive them. Most obvious example first: you can do postprocessing filters using Compute Shaders too. Like in PS, your target has a natural cartesian 2D parametrization. Unlike PS, you don't get implicit derivative calculations, so if you want anisotropic filtering (and there's good reasons to use aniso even for 2D processing, e.g. when your input and output buffers have different pixel aspect ratios) you need to use `SampleGrad`. But it doesn't have to derivatives in pixel coordinates; the derivative vectors you pass to `SampleGrad` define a Jacobian matrix. Any sufficiently (locally) smooth mapping between two 2-dimensional spaces allows such a Jacobian to be computed, and anisotropic filtering does the "right" thing (in the sense of it being a typically good approximation of Heckbert's Elliptical Weighted Averaging filter); there is still an approximation involved in determining filter shape from the Jacobian, and some more approximation in locally linearizing the approximation to place sample values.

Anyway, there's more common cases than don't involve arbitrary smooth 2D maps; in particular, grid-based meshes have a natural 2D parametrization. For Vertex Shaders, you'd typically see grid meshes in e.g. height map-based rendering.

Another place where meshes with a natural 2D parametrization crop up is with Tessellation (which I haven't discussed yet), particularly on "quad" domains. If the patch control points form a grid, you get a natural parametrization in the Hull Shader (and hence also a pair of natural axes along which to compute derivatives or at least differences). Similarly, the Domain Shader for quad domains gets 2D U/V coordinates which is also a natural 2D parametrization. And of course you can play the same game in Geometry Shaders too.

Reply

- **Sin permalink**

Let me say thank you for writing this series. I find it very helpful to deepen my knowledge in graphics programming.

I'm wondering though, where does the number 1.25 misses/sample in the case of bilinear filtering come from ? Can you elaborate on that ?

- **fgiesen permalink**

This is just a really really coarse rule of thumb. Also, this is a *bandwidth* estimate (see below), which I should've made clearer in the post, not a statement about the absolute frequency of cache misses – that depends on the cache architecture.

Anyway, 1.25 texels/pixel. Short version: I'm gonna give you some rationale for why you would expect that number to be 1 per pixel or a bit more (in the bandwidth sense). I'm not gonna go into why 1.25; Honestly, I just got that number off people working on texture samplers and didn't ask. But let's start with why you would expect something in the vicinity of 1 texel/pixel or more.

The mipmap level chosen per pixel is picked such that the texel-to-pixel ratio is as close to 1:1 as possible. Now, suppose that you're filling the entire screen with a single quad using a repeating 512×512 RGBA8888 texture, at 1:1 scale. If your texture cache is anything smaller than 512x512x4 bytes = 1MB (they're way smaller than that), the texture will not fit entirely inside the cache.

Suppose for now that we access the entire texture before we get back to a location we've already referenced (I'll get back to that in a minute). In that case, we will reference at least 1 "new" (not in the cache) texel for every pixel on the screen, more if we're unlucky; and this has nothing to do with bilinear filtering per se, this will happen with *any* filter kernel

(provided it's the same filter footprint for every pixel). Why? Let's look at a partially-rendered image (crappy ASCII art inbound):

```
#####
#####
#####
#####
@.....
```

'#' = pixel we finished earlier, '@' = pixel we're currently working on, '.' = pixel we're gonna look at later.

Now I've already said that for now, we're assuming 1:1 texel:pixel mapping – so *this picture looks the same if we plot texels we've accessed while rendering these pixels*. Bilinear is a 2×2 -tap kernel, so we sample 4 positions. Out of those 4, 3 are positions we already accessed while rendering the adjacent pixels (which happened "recently"). The fourth (the bottom right of the 4 pixels) is "new" and hasn't been accessed since we last accessed the texture.

This is hand-wavey, but I want to give you a mental picture of the process. A slightly more rigorous view is just that we're streaming 1MB of data through, say, a 32k cache. With 1:1 pixel scaling and point sampling (no bilinear), we just access one texel per pixel, and that's a texel we haven't accessed "in a long time" so it's not in the cache. As we make the filter kernel bigger (and better), the number of texels accessed per pixel increases linearly, but the number of "new" texels accessed does not. We may access more texels earlier than before, but we're still accessing the same texels eventually (is it clear why? If not, I can give more examples). Now if you keep growing the filter footprint at some point this will stop working as your footprint becomes large enough not to fit in the texture cache by itself, but that's not a problem for usual texture sampling filters. :)

Anyway. I was assuming before that the texel data falls out of the cache before we reference it again, despite using a repeating texture. Is that really true? Well, it depends. You might get lucky in some cases, but it's pretty much true in general; textures tend to be big. You gain a lot from making your texture caches large enough that texels hang around between shading for most adjacent pixels when rasterizing. Then it plateaus pretty soon until the cache is large enough to contain entire the entire texture working set for a shader.

For the sake of argument, let's say that our texture is stored in linear (raster-scan order), and we render strictly top to bottom, left to right. In that case, we might not even need the full texture cache! If we're rendering to a 2048×1024 screen (rounded from 1920×1080 for convenience), we get 4 repeats of the 512×512 texture per scan line, and it only takes $(512 \times 4) \times 2 = 4\text{kb}$ to store the 2 scan lines worth of texture that are actually accessed while rendering this scan line! So we get 4x reuse from the texture cache after all!

Well, not quite. Because now I'm gonna be mean and start slanting the texture coordinates (even when I perform a slant, it's still a 1:1 mapping in terms of density – the determinant of a slant is 1, and the amount of memory accessed should not change, just the order). Enough so that, by the point we've moved 512 pixels to the right, we've also moved 128 texel rows down in the texture. And now the 4 copies of the texture access texels distant enough from each other that (with our raster-order rasterization pattern) they're only gonna be in the cache by the next time we come around if the cache is large enough to contain the whole texture.

Actual GPUs do not rasterize in scan-line order, and they do not (by default) store textures in raster order; they use traversal patterns and layouts optimized for 2D access, to make sure that all orientations work roughly equally well. These make the worst cases better and the best cases worse. That's good – predictable latency you can easily design around, covering high variability is trickier since it leads to bursty behavior.

This covers 1:1 mappings. You can't get arbitrarily far from there – mip map selection will pick something reasonable, assuming that is that there are actually mip maps present and in use. But yes, if you were to zoom in, you'd access fewer unique texels per pixel; if you were to zoom out, you'd access more. At which point it's more handwaving – in general, you have about as many texels overshooting their target mip level than undershooting it, and so forth. But it gets pretty spotty at that point.

Okay, that gives us a very rough intuition for why 1 texel/pixel is reasonable. Why more?

Well; triangles have edges; caches have set conflicts; real scenes have tons of textures and several batches; and so forth. What this all presumes is a long steady state where we're merrily shading a large triangle, and the only real limit is cache capacity. So 1 is definitely optimistic. Why use an extra fudge factor of 0.25 instead of something else? I don't know, and I'm not gonna pretend I do. :)

Extra caveats: a) as said before, this is just a rule of thumb, b) it's also only valid for large textures and in the statistical sense (i.e. as average over longer periods of time), c) this is a statement about texture memory access *bandwidth* (throughput) not *latency*. What I've been saying here assumes that texels are cached individually. That is, of course, not the case in practice; texel caches are organized in terms of lines, same as other caches. Say we're sticking with the 4byte/pixel RGBA8888 format, and we have 64-byte cache lines. That means we have 16 texels per cache line. If the memory layout is good, we will thus get 16 texels per cache miss, and these are actually 16 texels that we will all want (sooner or later). Thus, we have only 1/16th the number of texture cache misses in terms of "how many texture accesses try to read from a line that's not currently in the cache". But the overall amount of memory we access is still the same: 1/16th the number of misses that now fetch 16x the data. And if the texture layout is *not* good (or doesn't match our access patterns), we might fetch 64 bytes worth of texels but only actually use 16 bytes from that cache lines, because the other 48 bytes are for texels we don't care about – this is always a risk when making cache lines bigger.

Okay, now suppose that we have 64-byte cache lines (with 16 texels/line at RGBA8888), but our texture sampler also processes 16 bilinear sampling requests per cycle (FWIW, both of these values are, as of this writing, realistic). Now, even assuming good texture caching, we're back at 1 cache miss/cycle (or worse); we're processing 16 pixels, with our larger cache lines roughly 1/16th of all pixels processed should miss the cache, and so we'd expect about 1 missed cache line per cycle.

Note that I can make the figure of "how many cache misses per cycle" go up and down by tweaking these values: halve the number of pixels processed per cycle – tada, half the cache misses! Double the cache line size – halved again! But note that to render our screen-filling quad, we still have the same number of texels accessed and read roughly the same amount of memory.

That's why the bandwidth number (number of memory read per pixel) is more useful as a general reference than any per-clock number; it's not as sensitive to architectural details, and it's also more relevant to one of the actual bottlenecks in GPUs (namely, memory bandwidth).

- o **Sin permalink**

Thank you for explaining it in great detail, and the caveats as well. I have to say that I somehow thought that the texels were cached individually :p

6. SteveM permalink

Fabian, thanks for this awesome series of articles. I was just referred back to this by a coworker after many years, and it sparked a pretty in-depth discussion.

I don't know if you'd still update this article, but it would've saved us a few hours of debate :)

You say:

"But this means that, unless you happen to hit the exact same location in a texture again and again, each texture sampling operation will miss about 1 texel on average – the actual measured value with bilinear filtering is around 1.25 misses/request (if you track pixels individually)."

Only in the comments section you correct:

"Also, this is a bandwidth estimate (see below), which I should've made clearer in the post, not a statement about the absolute frequency of cache misses"

You should make it explicit in the original article that the 1 texel/request figure is actually just a lower bound on the *bandwidth* requirements, and remove the reference to "cache misses", or clarify it. The 1.25 figure is an empirical estimate, at best, and, as you point out, you have no solid theoretical justification to back it up. So it should be made explicit that it's an empirically sourced number. The theoretical lower bound on bandwidth is 1 texel/request for point or bilinear filtering.

You can construct scenarios where the cache miss frequency is much higher or lower than 1 miss per request (usually much lower with tiling/swizzling), and you can construct pathological scenarios with linear memory layouts at 90 degrees to each other where the actual bandwidth requirements are multiples of this lower bound due to throwing away most of the data in cache lines you fetch and have to re-fetch (as you also point out in the comments), but you can't ever get *below* this lower bound on bandwidth, as long as you maintain close to 1 texel per pixel density.

I don't think my prose is very concise or intelligible as is, but I'm sure you could make the original article more intelligible, and maybe save future readers a bit of headache, as this still is a de facto reference several years on! :)

Keep up the good work!

Steve

Reply

o [fgiesen permalink](#)

Nowhere do I say that 1.25 is a lower bound. It's not. I'm talking about texture-sampling induced memory bandwidth (memory accesses) not number of texel requests. The lower bound is 0 (not quite, but you can get arbitrarily close) and occurs when you have a tiny texture (let's take a trivial example, a 2x2 texture) that fits fully within the sampler L1 caches. Consider the case where you're drawing a million batches with that one texture. If the texture fits inside the sampler caches, there will be no more requests all the way to memory after the cache lines containing the data were initially loaded, no matter how much you draw with it. Hence the limit of memory requests per texel sampled goes to zero (since you can increase the number of texels sampled arbitrarily without increasing the number of cache misses).

At the opposite extreme, the upper bound for texel misses per bilinear texture sample is 4 per pixel evaluated. A texture fetch/sample is fundamentally a "gather" style operation, and knowing sufficient details about the caches and texture tiling patterns, you can construct a worst-case pattern where *not a single texel fetched* is ever inside the L1 or L2 caches at the time when it's requested. A bilinear fetch grabs 4 texels and you can in fact arrange for none of them to ever hit the cache.

In practice, your texture is typically way too large to fit inside your sampler caches, and texture memory accesses have reasonable locality of reference and temporal locality. Under these conditions and with mip mapping enabled the number of texel misses per pixel will usually be close to the 1.25 I state (though this is really just a rule of thumb). Without mip mapping (or with a large negative mip map LOD bias) it's actually really easy to hit the worst case of 4 texels per pixel.

Reply

Trackbacks & Pingbacks

1. [A trip through the Graphics Pipeline 2011: Index](#) « The ryg blog
2. [A trip through the Graphics Pipeline 2011, part 8](#) « The ryg blog
3. [\(Updated\) 3D Graphics Pipeline Explained - 3D Tech News, Pixel Hacking, Data Visualization and 3D Programming - Geeks3D.com](#)

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 5

July 5, 2011

This post is part of the series “A trip through the Graphics Pipeline 2011” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

After the last post about texture samplers, we’re now back in the 3D frontend. We’re done with vertex shading, so now we can start actually rendering stuff, right? Well, not quite. You see, there’s a bunch still left to do before we actually start rasterizing primitives. So much so in fact that we’re not going to see any rasterization in this post – that’ll have to wait until next time.

Primitive Assembly

When we left the vertex pipeline, we had just gotten a block of shaded vertices back from the shader units, with the implicit promise that this block contains an integral number of primitives – i.e., we don’t allow triangles, lines or patches to be split across multiple blocks. This is important, because it means we can truly treat each block independently and never need to buffer more than one block of shader output – we can, of course, but we don’t have to.

The next step is to assemble all the vertices belonging to a single primitive (hence “primitive assembly”). If that primitive happens to be a point, this just reads exactly one vertex and passes it on. If it’s lines, it reads two vertices. If it’s triangles, three. And so on for patches with larger numbers of control points.

In short, all that happens here is that we gather vertices. We can either do this by reading the original index buffer and keeping a copy of our vertex index->cache position map around (as I described), or we can store the indices for the fully expanded primitives along with the shaded vertex data, which might take a bit more space for the output buffer but means we don’t have to read the indices again here. Either way works fine.

And now we have expanded out all the vertices that make up a primitive. In other words, we now have complete triangles, not just a bunch of vertices. So can we rasterize them already? Not quite.

Viewport culling and clipping

Oh yeah, that. Yeah, I guess we’d better do that first, huh? This is one part of pipeline that really does exactly what you’d expect, pretty much the way you would expect it too (i.e. the way it’s explained in the docs). So I’m not gonna explain polygon clipping in general here, you can look that up in any computer graphics textbook, although most make a terrible mess of it; if you want a good explanation, use Jim Blinn’s (chapter 13 of [this book](http://www.amazon.com/Jim-Blinns-Corner-Graphics-Pipeline/dp/1558603875) (<http://www.amazon.com/Jim-Blinns-Corner-Graphics-Pipeline/dp/1558603875>)), although you probably want to pass on his alternative [0,w] clip space these days, to avoid confusion if nothing else.

Anyway, clipping. The short version is this: Your vertex shader returns vertex positions on homogeneous clip space. Clip space is chosen to make the equations that describe the view frustum as simple as possible; in the case of D3D, they are $-w \leq x \leq w$, $-w \leq y \leq w$, $0 \leq z \leq w$, and $0 < w$; note that all the last equation really does is exclude the homogeneous point (0,0,0,0), which is something of a degenerate case.

We first need to find out if the triangle is partially or even completely outside any of these clip planes. This can be done very efficiently using **Cohen-Sutherland** (<http://en.wikipedia.org/wiki/Cohen-Sutherland>)-style out-codes. You compute the clip out-code (or just clip-code) for each vertex (this can be done at vertex shading time and stored along with the positions, for example). Then, for each primitive, the bitwise AND of the clip-codes will tell you all the view-frustum planes that *all* vertices in the primitive are on the wrong side of (if there’s any, that means the primitive is completely outside the view frustum and can be thrown away), and the bitwise OR of the clip-codes will tell you the planes that you need to clip the primitive against. Given the clipcodes, all this is just a few gates worth of hardware – simple stuff.

Additionally, the shaders can also generate a set of “cull distances” (a triangle will be discarded if any one cull distance for all vertices is less than zero), and a set of “clip distances” (which define additional clipping planes). These get considered for primitive rejection/clip testing too.

The actual clipping process, if invoked, can take one of two forms: we can either use an actual polygon clipping algorithm (which adds extra vertices and triangles), or we can add the clipping planes as extra edge equations to the rasterizer (if that sounds like gibberish to you, wait until the next part where I explain rasterization – it’ll make sense eventually). The latter is more elegant and doesn’t require an actual polygon clipper at all, but we need to be able to handle all normalized 32-bit floating point values as valid vertex coordinates; there might be a trick for building a fast HW rasterizer that does this, but it seems tricky to say the least. So I’m assuming there’s an actual clipper, with all that involves (generation of extra triangles etc). This is a nuisance, but it’s also very infrequent (more so than you think, I’ll get to that in a second), so it’s not a big deal. Not sure if that’s special hardware either, or if that path grabs a shader unit to do the actual clipping; depends on whether dispatching a new vertex shading load at this stage is awkward or not, how big a dedicated clipping unit is, and how many of them you need. I don’t know the answer to these questions, but at least from the performance side of things, it doesn’t much matter: we don’t really clip that often. That’s because we can use guard-band clipping.

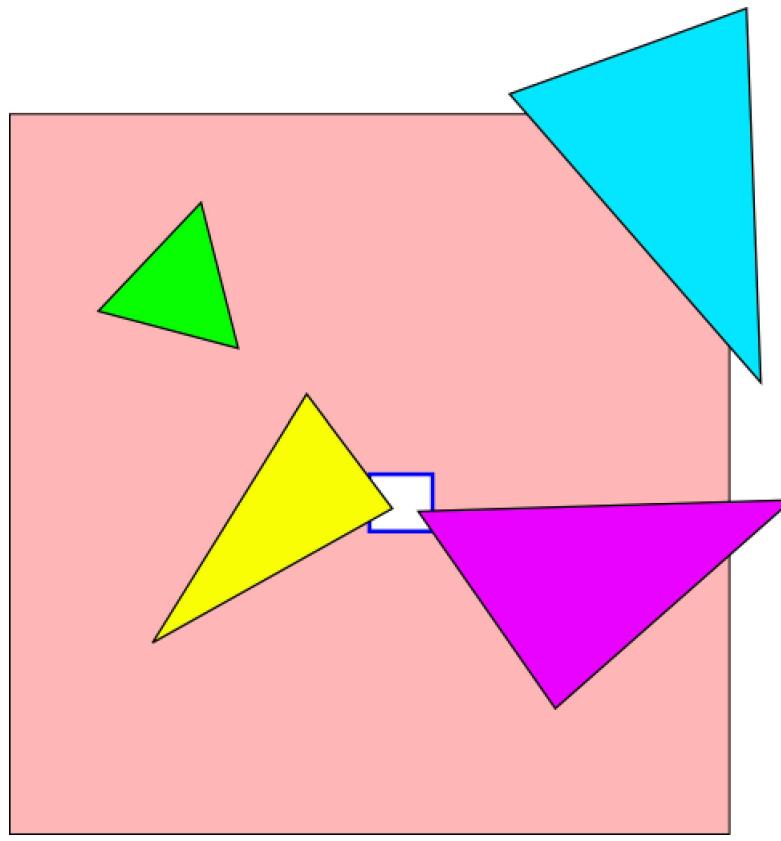
Guard-band clipping

The name is something of a misnomer; it’s not a fancy way of doing clipping. In fact, it’s quite the opposite: a straight-forward way of not doing clipping. :)

The underlying idea is very simple: Most primitives that are partially outside the left, right, top and bottom clip planes don’t need to be clipped at all. Triangle rasterization on GPUs works by, in effect, scanning over the full screen area (or more precisely, the scissor rect) and asking for every pixel: “is this pixel covered by the current triangle?” (In reality it’s a bit more complicated and way more efficient than that, but that’s the general idea). And that works just as well for triangles completely within the viewport as it does for triangles that extend past, say, the right and top clipping planes. As long as our triangle coverage test is reliable, we don’t need to clip against the left, right, top and bottom planes at all!

That test is usually done in integer arithmetic with some fixed precision. And eventually, as you move say one triangle vertex further and further out, you’ll get integer overflows and wrong test results. I think we can all agree that the rasterizer producing pixels that aren’t actually inside the triangle is, at the very least, extremely offensive behavior and should be illegal! Which it in fact is – hardware that does this is in violation of the spec.

There’s two solutions for this problem: The first is to make sure that your triangle tests never, ever generate the wrong results, no matter how your input triangle looks. If you manage that, then you don’t ever need to clip against the aforementioned four planes. This is called “infinite guard-band” because, well, the guard-band is effectively infinite. Solution two is to clip triangles eventually, just as they’re about to go outside the safe range where the rasterizer calculations can’t overflow. For example, say that your rasterizer has enough internal bits to deal with integer triangle coordinates that have $-32768 \leq X \leq 32767, -32768 \leq Y \leq 32767$ (note I’m using capital X and Y to denote screen-space positions; I’ll stick with this convention). You still do your viewport cull test (i.e. “is this triangle outside the view frustum”) with the regular view planes, but only actually clip against the guard-band clip planes which are chosen so that after the projection and viewport transforms, the resulting coordinates are in the safe range. I guess it’s time for an image:



(https://fgiesen.files.wordpress.com/2011/07/guardband_clip.png)
Guard-band clipping

The small white rectangle with blue outline that's roughly in the middle represents our viewport, while the big salmon-colored area around it is our guard band. It looks like a small viewport in this image, but I actually picked a huge one so you can see anything! With our $-32768 .. 32767$ guard-band clip range, that viewport would be about 5500 pixels wide – yes, that's some huge triangles right there :). Anyway, the triangles show off some of the important cases. The yellow triangle is the most common case – a triangle that extends outside the viewport but not the guard band. This just gets passed straight through, no further processing necessary. The green triangle is fully within the guard band, but outside the viewport region, so it would never get here – it's been rejected above by the viewport cull. The blue triangle extends outside the guard-band clip region and would need to be clipped, but again it's fully outside the viewport region and gets rejected by the viewport cull. Finally, the purple triangle extends both inside the viewport and outside the guard band, and so actually needs to be clipped.

As you can see, the kinds of triangles you need to actually have to clip against the four side planes are pretty extreme. As said, it's infrequent – don't worry about it.

Aside: Getting clipping right

None of this should be terribly surprising; nor should it sound too difficult, at least if you're familiar with the algorithms. But the devil's in the details, always. Here's some of the non-obvious rules the triangle clipper has to obey in practice. If it ever breaks *any* of these rules, there's cases where it will produce cracks between adjacent triangles that share an edge. This isn't allowed.

Vertex positions that are inside the view frustum must be preserved, bit-exact, by the clipper.

Clipping an edge AB against a plane must produce the same results, bit-exact, as clipping the edge BA (orientation reversed) against that plane. (This can be ensured by either making the math completely symmetric, or always clipping an edge in the same direction, say from the outside in).

Primitives that are clipped against multiple planes must always clip against planes in the same order. (Either that or clip against all planes at once)

If you use a guard band, you *must* clip against the guard band planes; you can't use a guard band for some triangles but then clip against the original viewport planes if you actually need to clip. Again, failing to do this will cause cracks – and if I remember correctly there was actually a piece of graphics hardware in the bad old days that shipped with this bug enshrined in silicon. Oops. :)

Those pesky near and far planes

Okay, so we have a really nice quick solution for the 4 side planes, but what about near and far? Particularly the near plane is bothersome, since with all the stuff that's only slightly outside the viewport handled, that's the plane we do most of our clipping for. So what can we do? A z guard band? But how would that work – we're not actually rasterizing along the z axis at all! In fact, it's just some value we interpolate over the triangle, damn!

On the plus side, though, it's just some value we interpolate over the triangle. And in fact the z-near test ($Z < 0$) is *really easy* to do once you interpolate Z – it's just the sign bit. z-far ($Z > 1$) is an extra compare though (not I'm using Z not z here, i.e. these are "screen" or post-projection coordinates). But still, we're doing Z-comparisons per pixel anyway (Z test!), so it's not a big extra expense. It depends, but doing z-clip this way is definitely an option. And you need to be able to skip z-near/z-far clipping if you want to support things like NVidias 'depth clamp' OpenGL extension; in fact, I would argue the existence of that extension is a pretty good hint that they're doing this, or at least used to for a while.

So we're down to one of the regular clip planes: $0 < w$. Can we get rid of this one too? The answer is yes, with a rasterization algorithm that works in homogeneous coordinates, e.g. [this one](http://www.cs.unc.edu/~olano/papers/2dh-tri/) (<http://www.cs.unc.edu/~olano/papers/2dh-tri/>). I'm not sure whether hardware uses that one though. It's nice an elegant, but it seems like it would be hard to obey the (very strict!) D3D11 rasterization rules to the letter using that algorithm. But maybe there's some cool tricks that I'm not aware of. Anyway, that's about it with clipping.

Projection and viewport transform

Projection just takes the x, y and z coordinates and divides them by w (unless you're using a homogeneous rasterizer which doesn't actually project – but I'll ignore that possibility in the following). This gives us normalized device coordinates, or NDCs, between -1 and 1. We then apply the viewport transform which maps the projected x and y to pixel coordinates (which I'll call X and Y) and the projected z into the range [0,1] (I'll call this value Z), such that at the z-near plane Z=0 and at the z-far plane Z=1.

At this point, we also snap pixels to fractional coordinates on the sub-pixel grid. As of D3D11, hardware is required to have exactly 8 bits of subpixel precision for triangle coordinates. This snapping turns some *very* thin slivers (which would otherwise cause problems) into degenerate triangles (which don't need to be rendered at all).

Back-face and other triangle culling

Once we have X and Y for all vertices, we can calculate the signed triangle area using a cross product of the edge vectors. If the area is negative, the triangle is wound counter-clockwise (here, negative areas correspond to counter-clockwise because we're now in the pixel coordinate space, and in D3D pixel space y increases downwards not upwards, so signs are inverted). If the area is positive, it's wound clockwise. If it's zero, it's degenerate and doesn't cover any pixels, so it can be safely culled. At this point, we know the triangle orientation so we can do back-face culling (if enabled).

And that's it! We're now ready for rasterization... almost. Actually we have to do triangle setup first. But doing that requires some knowledge of how rasterization will be performed, so I'll put that off until the next part... see you then!

Final remarks

Again, I skipped some parts and simplified others, so here's the usual reminder that things are a bit more complicated in reality: For example, I pretended that you just use the regular homogeneous clipping algorithm. Mostly, you do – but you can have some vertex shader attributes flagged as using screen-space linear instead of perspective-correct interpolation. Now, the regular homogeneous clip always does perspective-correct interpolation; in the case of screen-space linear attributes, you actually need to do some extra work to make it not perspective-correct. :)

I talk about primitives some of the time, but mostly I'm just focusing on triangles here. Points and lines aren't hard, but let's be honest, they're not what we're here for either. You can work out the details if you're interested. :)

There's tons of rasterization algorithms out there, some of which (like Olanos 2DH method that I cited) allow you to skip nearly all clipping, but as I mentioned, D3D11 has very strict requirements on the triangle rasterizer so there's not much wiggle room for HW implementations; I'm not sure if those methods can be tweaked to exactly follow the spec (there's a lot of subtle points that I'll cover next time). So here and in the following I'm assuming you can't do the ultra-sleek thing; then again, the not-quite-so-sleek approaches I'm running with have slightly less math per pixel in the rasterizer, so they might win for HW implementations anyway. And of course I might be missing the magic pixie dust right around the corner that solves all of these problems. That occurs surprisingly often in graphics. If you know an awesome solution, give me a shout in the comments!

Lastly, the triangle culling I'm describing here is the bare minimum; for example, the class of triangles that will generate zero pixels upon rasterization is much larger than just zero-area tris, and if you can find it out quickly enough (or with few enough gates), you can drop the triangle immediately and don't need to go through triangle setup. This is the last point where you can cull cheaply before going through triangle setup and at least some rasterization – finding other ways to early-reject tris pays off handsomely here.

From → Coding, Graphics Pipeline

17 Comments

1. Ben permalink

Ahh.. those 4 clipper requirements warm my heart!

I stumbled on those 4 exactly (except the first one, which was obvious before I started) when writing a rasterizer a couple years back.

Anyway – great series. Learning a lot.

Thanks

Reply

2. Jocelyn Houle permalink

"Again, failing to do this will cause cracks – and if I remember correctly there was actually a piece of graphics hardware in the bad old days that shipped with this bug enshrined in silicon. Oops. :)"

Mmhhh... Tasty bit of information... What GPU was that?

"If you know an awesome solution, give me a shout in the comments!"

Dachsbacher et al. wrote a 3D rasterizer (http://www.vis.uni-stuttgart.de/~dachsbcn/download/3dr_techreport.pdf) that failed to have academia's approval (i.e. only a tech report). But it does sport some interesting concepts that allows to rasterize triangles over arbitrary surfaces, and therefore bridges the gap between rasterization and ray tracing. Haven't heard of any IHV trying to accelerate such path, though... (kind of out of the loop, now) If only for robust ray-casting, this article is quite a gem...

Reply

- o [fgiesen](#) permalink

"What GPU was that?"

Not naming any names – not because I don't want to offend, I just seriously don't remember :). This was maybe 2000/2001 or so. When 3dfx, S3, Matrox and 3Dlabs were still around in the 3D graphics space, and you really would run into driver/compatibility issues even with totally basic stuff on a regular basis.

The 3D rasterizer paper looks sweet (as indeed most of Carstens papers do), but same as with Olanos paper, I don't see any obvious way to integrate things like subpixel grid snapping (required for D3D10/11 hardware!) into that kind of approach, especially since it needs to be absolutely, unconditionally robust. Another worry is precision – 32-bit float is definitely not enough to evaluate edge equations with full precision for large triangles; even doubles with their 53-bit mantissa (including hidden 1 bit) are uncomfortably tight. Plus any floating-point implementation needs to recompute the edge equations regularly using a multiply-accumulate style operation to avoid drift; in contrast, a fixed-point rasterizer only ever does integer adds after the initial triangle setup.

Reply

3. Kevin permalink

Very informative and detailed article :)

I think you may have forgotten one thing, though. All polygons must be clipped to the near plane by generating new vertices, not by discarding pixels, otherwise you'll run into accuracy problems when one of a triangle's vertices passes behind the view point(which is when it's behind the w=0 plane, I believe).

Reply

- o [fgiesen permalink](#)

I actually link to Olano's paper that describes a fully clipless rasterizer. You do *not* need to clip, but you do need to add an extra edge function corresponding to the near plane.

Nvidia used completely clipless rasterizers for quite a while. I'm not sure whether they still do. The D3D11 spec insists on some invariance properties that are quite tricky to obtain in a clipless rasterizer.

Reply

4. Sin permalink

"As of D3D11, hardware is required to have exactly 8 bits of subpixel precision for triangle coordinates."

I am not familiar with subpixel precision. I assume it means that there are several subpixels (in case of 8 bits, 256) in a pixel, and is totally unrelated to sub-pixel resolution (<https://www.grc.com/ctwhat.htm>) ?

But if that's the case, then what will happen if a triangle occupies some subpixels ? Will the pixel get colored using a certain weight, like with the ratio of occupied subpixels divided by total subpixels in a pixel ?

Reply

- o [fgiesen permalink](#)

No. Sub-pixel accuracy just means that the triangle coordinates aren't snapped to integer pixel positions (which produces very visible artifacts especially with slow movement). 8-bit subpixel precision means that coordinates are rounded to the next 1/256ths of a pixel in both the x and y directions. This has nothing to do with coverage computations or any kind of anti-aliasing (or sub-pixel rendering); I explain a bit more here.

Reply

- o [Sin permalink](#)

Thanks, I read your post and left another comment there.

I thought a bit after reading your post last time, and I did think that it perhaps had something to do with coverage :p

So, just to clear out my confusion, this subpixel precision only means that the coordinates aren't snapped to integer positions, which can further be used for things like alpha to coverage and MSAA, for example the technique explained in <http://software.intel.com/en-us/articles/rendering-grass-with-instancing-in-directx-10>

Now that I think about it, these features do need non-integer positions. Otherwise, how can they know that a certain triangle occupies a number of samples out of 16 samples (MSAA4x) in a certain pixel.

- o [Sin permalink](#)

err, I meant MSAA 16x. Something wrong with my head.

5. zhebin permalink

Great articles, and help me a lot!

But I have two questions in this topic:

1.

How to handle clip space coordinate (0, 0, 0, 0) in the Cohen-Sutherland algorithm?

Considering the "w > 0" plane, The outcode of (0, 0, 0, 0) will be 1xxxx, but the other two vertices's outcodes may be 0xxxx, so it still needs clipping. However, after clipping, I think we still get the same three vertices – I assume clipping for open interval domain($w > 0$) works the same way as closed interval($w \geq 0$).

Another idea is that, as you mentioned, this is a degenerated triangle(a line), so we can throw away it. But this is beyond the power of Cohen-Sutherland outcodes and needs special handling.

Unfortunately, I didn't see such kind of special handling in many open source software rasterizers.

2.

"If you use a guard band, you must clip against the guard band planes; you can't use a guard band for some triangles but then clip against the original viewport planes if you actually need to clip. Again, failing to do this will cause cracks."

Could you explain why clipping to original viewport planes will cause cracks?

It's difficult for me to imagine a case where cracks will happen, and AFAIK, there are still some software rasterizers doing things like this.

Reply

- o [fgiesen permalink](#)

First question:

Cohen-Sutherland can only clip for inclusive bounds. If you're gonna do a w-clip with Cohen-Sutherland, you end up having to pick an epsilon and clip to " $\text{eps} \leq w$ " instead. That works, but we can do better by looking at what happens with 0. The only clip space "point" the $w > 0$ rule actually catches is $(0,0,0,0)$ —using quotes here because that's not technically a point as per the definition of homogeneous points (all-0 is explicitly excluded), but we'll run with it. And we can just look at cases involving that point directly, rather than trying to catch them with the clipper. Let's do triangles: a triangle is a convex polygon with 3 vertices and every point inside the triangle can be written as a convex combination of those three vertices

$$\lambda_0 v_0 + \lambda_1 v_1 + \lambda_2 v_2$$

where all the $\lambda_i \geq 0$ and $\lambda_0 + \lambda_1 + \lambda_2 = 1$. This works in a projective space with homogeneous coordinates as well, except we don't need the weights to sum to 1 anymore; homogeneous coordinates are only unique up to scale. So in a projective space, the equivalent of looking at convex combinations is looking at conical combinations – linear combinations where the weights are non-negative.

Now suppose one of the triangle vertices hits $(0,0,0,0)$, say v_2 . That means the last term of our conical combination is always zero and the whole expression reduces to

$$\lambda_0 v_0 + \lambda_1 v_1$$

with the $\lambda_i \geq 0$. That's a conical combination of two homogeneous vertices, which corresponds to a line. But a triangle that is actually a line is degenerate and doesn't generate any visible fragments according to our rasterization rules! So long story short, if one of your triangles' vertices is $(0,0,0,0)$ (again, technically not a point), that guarantees the triangle is degenerate and you can drop it, same as any other zero-area triangle. The same argument applies to lines: a line with one of the vertices at $(0,0,0,0)$ is actually degenerate and just a homogeneous point (whether that means it should produce visible fragments or not depends on your choice of line rasterization rules), and a homogeneous point at $(0,0,0,0)$ is not actually a valid point and gets dropped.

Second question:

Consider a quad split into two triangles, with the interior diagonal going from the bottom left to the top right point.

The top right and bottom right points are entirely within the viewport. The bottom left point is *just* outside the left edge of the viewport (but well within the guard band). The top left point is *way* off to the left, outside the guard band.

The bottom right triangle is fully within the guard band and can get rasterized directly. The top left triangle is clipped.

Now if you clip the bottom left vertex to the actual viewport while processing the clipped triangle, the clipped vertex may not end up on exactly the same sub-pixel position, due to round-off. If you have a shared edge between two triangles and you slightly nudge one of its vertices in one triangle but not the other, you get cracks.

There's other ways to avoid this, but *by far* the easiest is to make sure that edges containing the same two vertices will always get handled the same way. In this case, that means you have one set of clip planes (corresponding to the guard band boundaries) that you use for clipping, and you don't switch from the "looser" guard band planes to the "tighter" viewport planes once you discover you actually need to clip.

Reply

- **fgiesen permalink**

Just to make that clear: the point of the answer to the first question is that you can treat the $w > 0$ clip purely as a symbolical operation. It does not actually have to be implemented as a "proper" plane clip at all, you just need to check for it and reject the primitive if necessary.

- **zhebin permalink**

First question:

Thanks! I got your point.

And I came up with the approach to catch the $(0,0,0,0)$ case with outcodes right after posting this question. That is, $(0,0,0,0)$ is equivalent to the following conditions:

$-w \leq x \leq w$ (true)

$-w \leq y \leq w$ (true)

$0 \leq z \leq 0$ (false)

So if the outcode of one vertex is 0000001b (meaning is shown as follows), we can conclude it is $(0,0,0,0)$ and this triangle can be dropped directly. This method may be not so efficient, but it can be integrated well into the outcodes algorithm.
near | far | left | right | bottom | top | w

0 0 0 0 0 0 1

Second question:

Yes, this case opens my mind, and I think I need take care of such cases in my own rasterizer implementation. Thank you very much!

- o **zhebin permalink**

Sorry that the typesetting is not very good.

Correct the typo:

" $-w \leq x \leq w$ " (true)

" $-w \leq y \leq w$ " (true)

" $0 \leq z \leq 0$ " (false)

6. xdnoam permalink

"The answer is yes, with a rasterization algorithm that works in homogeneous coordinates, e.g. this one."

The link is broken, here's an updated one – <https://www.csee.umbc.edu/~olano/papers/2dh-tri/2dh-tri.pdf>

Reply

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog

2. Triangle rasterization in practice « The ryg blog

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 7

July 8, 2011

This post is part of the series “A trip through the Graphics Pipeline 2011” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

In this installment, I'll be talking about the (early) Z pipeline and how it interacts with rasterization. Like the last part, the text won't proceed in actual pipeline order; again, I'll describe the underlying algorithms first, and then fill in the pipeline stages (in reverse order, because that's the easiest way to explain it) after the fact.

Interpolated values

Z is interpolated across the triangle, as are all the attributes output by the vertex shader. So let me take a minute to explain how that works. At this point I originally had a section on how the math behind interpolation is derived, and why perspective interpolation works the way it works. I struggled with that for hours, because I was trying to limit it to maybe one or two paragraphs (since it's an aside), and what I can say now is that if I want to explain it properly, I need more space than that, and at least one or two pictures; a picture may say more than thousand words, but a nice diagram takes me about as long to prepare as a thousand words of text, so that's not necessarily a win from my perspective :). Anyway, this is something of a tangent anyway, so I'm adding it to my pile of “graphics-related things to write up properly at some point”. For now, I'm giving you the executive summary:

Just linearly interpolating attributes (colors, texture coordinates etc.) across the screen-space triangle does not produce the right results (unless the interpolation mode is one of the “no perspective” ones, in which case ignore what I just wrote). However, say we want to interpolate a 2D texture coordinate pair (\bar{s}, \bar{t}) . It turns out you do get the right results if you linearly interpolate $\frac{1}{w}, \frac{\bar{s}}{w}$ and $\frac{\bar{t}}{w}$ in screen-space (w here is the homogeneous clip-space w from the vertex position), then per-pixel take the reciprocal of $\frac{1}{w}$ to get w, and finally multiply the other two interpolated fractions by w to get s and t. The actual linear interpolation boils down to setting up a plane equation and then plugging the screen-space coordinates in. And if you're writing a software perspective texture mapper, that's the end of it. But if you're interpolating more than two values, a better approach is to compute (using perspective interpolation) **barycentric coordinates** ([http://en.wikipedia.org/wiki/Barycentric_coordinate_system_\(mathematics\)](http://en.wikipedia.org/wiki/Barycentric_coordinate_system_(mathematics))) – let's call them λ_0 and λ_1 – for the current pixel in the original clip-space triangle, after which you can interpolate the actual vertex attributes using regular linear interpolation without having to multiply everything by w afterwards.

So how much work does that add to triangle setup? Setting up the $\frac{\lambda_0}{w}$ and $\frac{\lambda_1}{w}$ for the triangle requires 4 reciprocals, the triangle area (which we already computed for back-face culling!), and a few subtractions, multiplies and adds. Setting up the vertex attributes for interpolation is really cheap with the barycentric approach – two subtractions per attribute (if you don't use barycentric, you get some more multiply-add action here). Follow me? Probably not, unless you've implemented this before. Sorry about that – but it's fairly safe to ignore all this if you don't understand it.

Let's get back to why we're here: the one value we want to interpolate *right now* is Z, and because we computed Z as $\frac{z}{w}$ at the vertex level as part of projection (see previous part), so it's already divided by w and we can just interpolate it linearly in screen space. Nice. What we end up with is a plane equation for $Z = aX + bY + c$ that we can just plug X and Y into to get a value. So, here's the punchline of my furious hand-waving in the last few paragraphs: Interpolating Z at any given point boils down to two multiply-adds. (Starting to see why GPUs have fast multiply-accumulate units? This stuff is absolutely everywhere!).

Early Z/Stencil

Now, if you believe the place that graphics APIs traditionally put Z/Stencil processing into – right before alpha blend, way at the bottom of the pixel pipeline – you might be confused a bit. Why am I even discussing Z at the point in the pipeline where we are right now? We haven't even started shading pixels! The answer is simple: the Z and stencil tests reject pixels. Potentially the majority of them. You really, *really* don't want to completely shade a detailed mesh with complicated materials, to then throw away

95% of the work you just did because that mesh happens to be mostly hidden behind a wall. That's just a really stupid waste of bandwidth, processing power and energy. And in most cases, it's completely unnecessary: most shaders don't do anything that would influence the results of the Z test, or the values written back to the Z/stencil buffers.

So what GPUs actually do when they can is called "early Z" (as opposed to late Z, which is actually at the late stage in the pipeline that traditional API models generally display it at). This does exactly what it sounds like – execute the Z/stencil tests and writes early, right after the triangle has been rasterized, and before we start sending off pixels to the shaders. That way, we notice all the rejected pixels early, without wasting a lot of computation on them. However, we can't always do this: the pixel shader may ignore the interpolated depth value, and instead provide its own depth to be written to the Z-buffer (e.g. depth sprites); or it might use discard, alpha test, or alpha-to-coverage, all of which "kill" pixels/samples during pixel shader execution and mean that we can't update the Z-buffer or stencil buffer early because we might be updating depth values for samples that later get discarded in the shader!

So GPUs actually have two copies of the Z/stencil logic; one right after the rasterizer and in front of the pixel shader (which does early Z) and one after the shader (which does late Z). Note that we can still, in principle, do the depth testing in the early-Z stage even if the shader uses some of the sample-killing mechanism. It's only writes that we have to be careful with. The only case that really precludes us from doing any early Z-testing at all is when we write the output depth in the pixel shader – in that case the early Z unit simply has nothing to work with.

Traditionally, APIs just pretended none of this early-out logic existed; Z/Stencil was in a late stage in the original API model, and any optimizations such as early-Z had to be done in a way that was 100% functionally consistent with that model; i.e. drivers had to detect when early-Z was applicable, and could only turn it on when there were no observable differences. By now APIs have closed that gap; as of DX11, shaders can be declared as "force early-Z", which means they run with full early-Z processing even when the shader uses primitives that aren't necessarily "safe" for early-Z, and shaders that write depth can declare that the interpolated Z value is conservative (i.e. early Z reject can still happen).

Z/stencil writes: the full truth

Okay, wait. As I've described it, we now have two parts in the pipeline – early Z and late Z – that can both write to the Z/stencil buffers. For any given shader/render state combination that we look at, this will work – in the steady state. But that's not how it works in practice. What actually happens is that we render a few hundred to a few thousand batches per frame, switching shaders and render state regularly. Most of these shaders will allow early Z, but some won't. Switching from a shader that does early Z to one that does late Z is no problem. But going back from late Z to early Z is, if early Z does any writes: early Z is, well, earlier in the pipeline than late Z – that's the whole point! So we may start early-Z processing for one shader, merrily writing to the depth buffer while there's still stuff down in the pipeline for our old shader that's running late-Z and may be trying to write the same location at the same time – classic race condition. So how do we fix this? There's a bunch of options:

Once you go from early-Z to late-Z processing within a frame (or at least a sequence of operations for the same render target), you stay at late-Z until the next point where you flush the pipeline anyway. This works but potentially wastes lots of shader cycles while early-Z is unnecessarily off.

Trigger a (pixel) pipeline flush when going from a late-Z shader to an early-Z shader – also works, also not exactly subtle. This time, we don't waste shader cycles (or memory bandwidth) but stall instead – not much of an improvement.

But in practice, having Z-writes in two places is just bad news. Another option is to not ever write Z in the early-Z phase; always do the Z-writes in late-Z. Note that you need to be careful to make conservative Z-testing decisions during early Z if you do this! This avoids the race condition but means the early Z-test results may be stale because the Z-write for the currently-dispatched pixels won't happen until a while later.

Use a separate unit that keeps track of Z-writes for us and enforces the correct ordering; both early-Z and late-Z must go through this unit.

All of these methods work, and all have their own advantages and drawbacks. Again I'm not sure what current hardware does in these cases, but I have strong reason to believe that it's one of the last two options. In particular, we'll meet a functional unit later down the road (and the pipeline) that would be a good place to implement the last option.

But we're still doing all this testing per pixel. Can't we do better?

Hierarchical Z/Stencil

The idea here is that we can use our tile trick from rasterization again, and try to Z-reject whole tiles at a time, before we even descend down to the pixel level! What we do here is a strictly conservative test; it may tell us that “there might be pixels that pass the Z/stencil-test in this tile” when there are none, but it will never claim that all pixels are rejected when in fact they weren’t.

Assume here that we’re using “less”, “less-equal”, or “equal” as Z-compare mode. Then we need to store the *maximum* Z-value we’ve written for that tile, per tile. When rasterizing a triangle, we calculate the *minimum* Z-value the active triangle is going to write to the current tile (one easy conservative approximation is to take the min of the interpolated Z-values at the four corners of the current tile). If our triangle minimum-Z is larger than the stored maximum-Z for the current tile, the triangle is guaranteed to be completely occluded. That means we now need to track maximum-Z per-tile, and keep that value up to date as we write new pixels – though again, it’s fine if that information isn’t completely up to date; since our Z-test is of the “less” variety, values in the Z buffer will only get smaller over time. If we use a per-tile maximum-Z that’s a bit out of date, it just means we’ll get slightly worse early rejection rates than we could; it doesn’t cause any other problems.

The same thing works (with min/max and compare directions swapped) if we’re using one of the “greater”, “greater-equal” or “equal” Z-tests. What we can’t easily do is change from one of the “less”-based tests to a “greater”-based tests in the middle of the frame, because that would make the information we’ve been tracking useless (for less-based tests we need maximum-Z per tile, for greater-based tests we need minimum-Z per tile). We’d need to loop over the whole depth buffer to recompute min/max for all tiles, but what GPUs actually do is turn hierarchical-Z off once you do this (up until the next Clear). So: don’t do that.

Similar to the hierarchical-Z logic I’ve described, current GPUs also have hierarchical stencil processing. However, unlike hierarchical-Z, I haven’t seen much in the way of published literature on the subject (meaning, I haven’t run into it – there might be papers on it, but I’m not aware of them); as a game console developer you get access to low-level GPU docs which include a description of the underlying algorithms, but frankly, I’m definitely not comfortable writing about something here where really the only good sources I have are various GPU docs that came with a thick stack of NDAs. Instead I’ll just nebulously note that there’s magic pixie dust that can do certain kinds of stencil testing very efficiently under controlled circumstances, and leave you to ponder what that might be and how it might work, in the unlikely case that you deeply care about this – presumably because your father was killed by a hierarchical stencil unit and you’re now collecting information on its weak points for your revenge, or something like that.

Putting it all together

Okay, we now have all the algorithms and theory we need – let’s see how we can take our new set of toys and wire it up with what we already have!

First off, we now need to do some extra triangle setup for Z/attribute interpolation. Not much to be done about it – more work for triangle setup; that’s how it goes. After that’s coarse rasterization, which I’ve discussed in the previous part.

Then there’s hierarchical Z (I’m assuming less-style comparisons here). We want to run this between coarse and fine rasterization. First, we need the logic to compute the minimum Z estimates for each tile. We also need to store the per-tile maximum Zs, which don’t need to be exact: we can shave bits as long as we always round up! As usual, there’s a trade-off here between space used and early-rejection efficiency. In theory, you could put the Z-max info into regular memory. In practice, I don’t think anyone does this, because you want to make the hierarchical-Z decision without a ton of extra latency. The other option is to put dedicated memory for hierarchical Z onto the chip – usually as SRAM, the kind of memory you also make caches out of. For 24-bit Z, you probably need something like 10-14 bits per tile to store a reasonable-accuracy Z-max in a compact encoding. Assuming 8x8 tiles, that means less than 1MBit (128k) of SRAM to support resolutions up to 2048x2048 – sounds like a plausible order of magnitude to me. Note that these things are fixed size and shared for the whole chip; if you do a context switch, you lose. If you allocate the wrong depth buffers to this memory, you can’t use hierarchical Z on the depth buffers that actually matter, and you lose. That’s just how it goes. This kind of things is why hardware vendors regularly tell you to create your most important render targets and depth buffers first; they have a limited supply of this type of memory (there’s more like it, as you’ll see), and when it runs out, you’re out of luck. Note they don’t necessarily need to do this all-or-nothing; for example, if you have a really large depth buffer, you might only get hierarchical Z in the top left 2048x1536 pixels, because that’s how much fits into the Z-max memory. It’s not ideal, but still much better than disabling hierarchical-Z outright.

And by the way, “Real-Time Rendering” mentions at this point that “it is likely that GPUs are using hierarchical Z-buffers with more than two levels”. I doubt this is true, for the same reason that I doubt they use a multilevel hierarchical rasterizer: adding more levels makes the easy cases (large triangles) even faster while adding latency and useless work for small triangles: if you’re drawing a triangle that fits inside a single 8×8 tile, any coarser hierarchy level is pure overhead, because even at the 8×8 level, you’d just do one test to trivial-reject the triangle (or not). And again, for hardware, it’s not that big a performance issue; as long as you’re not consuming extra bandwidth or other scarce resources, doing more compute work than strictly necessary isn’t a big problem, as long as it’s within reasonable limits,

Hierarchical stencil is also there and should also happen prior to fine rast, most likely in parallel with hierarchical Z. We’ve established that this runs on air, love and magic pixie dust, so it doesn’t need any actual hardware and is probably always exactly right in its predictions. Ahem. Moving on.

After that is fine rasterization, followed in turn by early Z. And for early Z, there’s two more important points I need to make.

Revenge of the API order

For the past few parts, I’ve been playing fast and loose with the order that primitives are submitted in. So far, it didn’t matter; not for vertex shading, nor primitive assembly, triangle setup or rasterization. But Z is different. For Z-compare modes like “less” or “lessequal”, it’s very important what order the pixels arrive in; if we mess with that, we risk changing the results and introducing nondeterministic behavior. More importantly, as per the spec, we’re free to execute operations in any order *so long as it isn’t visible to the app*; well, as I just said, for Z processing, order is important, so we need to make sure that triangles arrive at Z processing in the right order (this goes for both early and late Z).

What we do in cases like this is go back in the pipeline and look for a reasonable spot to sort things into order again. In our current path, the best candidate location seems to be primitive assembly; so when we start assembling primitives from shaded vertex blocks, we make sure to assemble them strictly in the original order as submitted by the app to the API. This means we might stall a bit more (if the PA buffer holds an output vertex block, but it’s not the correct one, we need to wait and can’t start setting up primitives yet), but that’s the price of correctness.

Memory bandwidth and Z compression

The second big point is that Z/Stencil is a serious bandwidth hog. This has a couple of reasons. For one, this is the one thing we really run for all samples generated by the rasterizer (assuming Z/Stencil isn’t off, of course). Shaders, blending etc. all benefit from the early rejection we do; but even Z-rejected pixels do a Z-buffer read first (unless they were killed by hierarchical Z). That’s just how it works. The other big reason is that, when multisampling is enabled, the Z/stencil buffer is per sample; so 4x MSAA means 4x the memory bandwidth cost of Z? For something that takes a substantial amount of memory bandwidth even at no MSAA, that’s seriously bad news.

So what GPUs do is Z compression. There’s various approaches, but the general idea is always the same: assuming reasonably-sized triangles, we expect a lot of tiles to just contain one or maybe two triangles. If that happens, then instead of storing Z-values for the whole tile, we just store the plane equation of the triangle that filled up this tile. That plane equation is (hopefully) smaller than the actual Z data. Without MSAA, one tile covers 8×8 actual pixels, so triangles need to be relatively big to cover a full tile; but with 4x MSAA, a tile effectively shrinks to 4×4 pixels, and covering full tiles gets easier. There’s also extensions that can support 2 triangles etc., but for reasonably-sized tiles, you can’t go much larger than 2-3 tris and still actually save bandwidth: the extra plane equations and coverage masks aren’t free!

Anyway, point is: this compression, when it works, is fully lossless, but it’s not applicable to all tiles. So we need some extra space to denote whether a tile is compressed or not. We could store this in regular memory, but that would mean we now need to wait two full memory round-trips latencies to do a Z-read. That’s bad. So again, we add some dedicated SRAM that allows us to store a few (1-3) bits per tile. At its simplest, it’s just a “compressed” or “not compressed” flag, but you can get fancy and add multiple compression modes and such. A nice side effect of Z-compression is that it allows us to do fast Z-clears: e.g. when clearing to Z=1, we just set all tiles to “compressed” and store the plane equation for a constant Z=1 triangle.

All of the Z-compression thing, much like texture compression in the texture samplers, can be folded into memory access/caching logic, and made completely transparent to everyone else. If you don’t want to send the plane equations (or add the interpolator logic) to the Z memory access block, it can just infer them from the Z data and use some integer delta-coding scheme. This kind of

approach usually needs extra bits per sample to actually allow lossless reconstruction, but it can lead to simpler data paths and nicer interface between units, which hardware guys love.

And that's it for today! Next up: Pixel shading and what happens around it.

Postscript

As I said earlier, the topic of setting up interpolated attributes would actually make for a nice article on its own. I'm skipping that for now – might decide to fill this gap later, who knows.

Z processing has been in the 3D pipeline for ages, and a serious bandwidth issue for most of the time; people have thought long and hard about this problem, and there's a zillion tricks that go into doing "production-quality" Z-buffering for GPUs, some big, some small. Again, I'm just scratching the surface here; I tried to limit myself to the bits that are useful to know for a graphics programmer. That's why I don't spend much time on the details of hierarchical Z computations or Z compression and the like; all of this is very specific on hardware details that change slightly in every generation, and ultimately, mostly there's just no practical way you get to exploit any of this usefully: If a given Z-compression scheme works well for your scene, that's some memory bandwidth you can spend on other things. If not, what are you gonna do? Change your geometry and camera position so that Z-compression is more efficient? Not very likely. To a hardware designer, these are all algorithms to be improved on in every generation, but to a programmer, they're just facts of life to deal with.

This time, I'm not going into much detail on how memory accesses work in this stage of the pipeline. That's intentional. There's a key to high-throughput pixel shading and other per-pixel or per-sample processing, but it's later in the pipeline, and we're not there yet. Everything will be revealed in due time :)

From → Coding, Graphics Pipeline

13 Comments

1. KeyJ permalink

Why is the order of the primitives so important and needs to be preserved? As long as only fully opaque triangles are rendered (i.e. blending is off, which should be the case for the largest part of typical scenes), the rendering order shouldn't make any difference (except maybe some additional overdraw if the application sorted primitives by depth), or am I mistaken here?

Reply

- o [fgiesen permalink](#)

Suppose you're rendering a batch that has 3 triangles, all with the exact same vertex coordinates but different colors (this example is a bit contrived, but nevertheless completely valid). Let's say the first triangle is red, the second green, and the third blue.

If you render this with "less" as compare mode onto a freshly cleared render target and depth buffer, you will only see the red triangle. If you use "less_equal", you will only see the blue triangle. That's what the functional API model (that is fully sequential) requires.

Allowing reordering of triangles in the pipeline, even just for opaque rendering, breaks this: for both compare modes, you can make any of the 3 triangles come out "on top" using an appropriate ordering.

This is a bit tricky, since you *can* allow triangle reordering within a batch if you only care about the contents of the Z-buffer (e.g. for shadow map rendering with a NULL pixel shader). For both compare functions, the Z-buffer at the end of the frame will contain at each location the minimum of all depth values that have been written there.

However, the output of our Z-processing isn't just the Z-buffer, it's also an updated coverage mask that tells us which pixels/samples to shade. And for the coverage mask computation, order matters.

Reply

2. TomF permalink

I believe it's common in HW to have all real Z done late. That result is then fed back "upstream" into the coarse Z unit, and the only sort of early Z is coarse Z. This is conservative and can be "late" – a triangle occluded by the previous triangle will still be shaded – but it is always safe and requires no mode-switching shenanigans.

Reply

- o [fgiesen permalink](#)

I'm not sure which type of implementation is more common these days, but I've definitely worked with chips that support both pre- and post-shading Z (including ones that don't do any coarse Z). But I've never seen anyone use anything but the obvious brute-force solution for the late Z->early Z transition: just flush the pixel pipeline when you're doing that switch.

Reply

3. piyush permalink

contrary to popular belief, computer scholars do have a sense of humor.
and you, sir are a living example

Reply

4. Martin Wardener permalink

"early Z is, well, earlier in the pipeline than early Z" .. I assume you mean "...than late Z" ..?

Reply

- o [fgiesen permalink](#)

Indeed! Thank you, fixed.

Reply

5. PixInverse permalink

I'm little confused by your explanation of perspective correct interpolation. Once we have performed the perspective divide $x/w, y/w, z/w$, then $Z(z/w)$, the value we want to interpolate, doesn't vary linearly anymore across the surface of the 2D triangle. What now varies linearly is $1/Z$.

To interpolate a vertex attribute correctly, we first need to divide by the vertex attribute value by Z of the vertex it is defined to, then linearly interpolate them, and then finally multiply the result by Z -Interpolated, which is the depth of the point on the triangle, that the pixel overlaps.

$$Z_{\text{correct}} = 1 / (\text{Interpolated } 1/Z)$$

$$U_{\text{incorrect}} = (\text{Interpolated } u/Z) * Z_{\text{correct}}$$

$$V_{\text{incorrect}} = (\text{Interpolated } v/Z) * Z_{\text{correct}}$$

In your post above, you say $u/w, v/w$ and z/w can be interpolated linearly for perspective correct interpolation. Could you please explain if I'm missing anything here?

Reply

- o [fgiesen permalink](#)

$Z=z/w$ ends up being an affine function of screen-space X and Y and is what's used for depth buffering.

It's never used for interpolation. Remember Z is set up by the projection matrix to reach (depending on the convention) either 0 and 1 or -1 and 1 at z_{near} and z_{far} , respectively.

What you use for interpolation is $1/w$ (which is also linear in screen space). Normally you set up perspective-corrected barycentric coordinates (usually called I and J) and then you interpolate $I/w, J/w$ and $1/w$ in screen space, and then for every pixel you solve $I=(I/w)/(1/w), J=(J/w)/(1/w)$ and use I and J to interpolate attributes. (As many as you want, and without having to individually divide them all through by the vertex w's).

Reply

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog
2. A trip through the Graphics Pipeline 2011, part 8 « The ryg blog
3. A trip through the Graphics Pipeline 2011, part 9 « The ryg blog
4. Linear Depth | The Devil In The Details

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 7

July 8, 2011

This post is part of the series “A trip through the Graphics Pipeline 2011” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

In this installment, I'll be talking about the (early) Z pipeline and how it interacts with rasterization. Like the last part, the text won't proceed in actual pipeline order; again, I'll describe the underlying algorithms first, and then fill in the pipeline stages (in reverse order, because that's the easiest way to explain it) after the fact.

Interpolated values

Z is interpolated across the triangle, as are all the attributes output by the vertex shader. So let me take a minute to explain how that works. At this point I originally had a section on how the math behind interpolation is derived, and why perspective interpolation works the way it works. I struggled with that for hours, because I was trying to limit it to maybe one or two paragraphs (since it's an aside), and what I can say now is that if I want to explain it properly, I need more space than that, and at least one or two pictures; a picture may say more than thousand words, but a nice diagram takes me about as long to prepare as a thousand words of text, so that's not necessarily a win from my perspective :). Anyway, this is something of a tangent anyway, so I'm adding it to my pile of “graphics-related things to write up properly at some point”. For now, I'm giving you the executive summary:

Just linearly interpolating attributes (colors, texture coordinates etc.) across the screen-space triangle does not produce the right results (unless the interpolation mode is one of the “no perspective” ones, in which case ignore what I just wrote). However, say we want to interpolate a 2D texture coordinate pair $(\frac{s}{w}, \frac{t}{w})$. It turns out you do get the right results if you linearly interpolate $\frac{1}{w}$, $\frac{s}{w}$ and $\frac{t}{w}$ in screen-space (w here is the homogeneous clip-space w from the vertex position), then per-pixel take the reciprocal of $\frac{1}{w}$ to get w, and finally multiply the other two interpolated fractions by w to get s and t. The actual linear interpolation boils down to setting up a plane equation and then plugging the screen-space coordinates in. And if you're writing a software perspective texture mapper, that's the end of it. But if you're interpolating more than two values, a better approach is to compute (using perspective interpolation) **barycentric coordinates** ([http://en.wikipedia.org/wiki/Barycentric_coordinate_system_\(mathematics\)](http://en.wikipedia.org/wiki/Barycentric_coordinate_system_(mathematics))) – let's call them λ_0 and λ_1 – for the current pixel in the original clip-space triangle, after which you can interpolate the actual vertex attributes using regular linear interpolation without having to multiply everything by w afterwards.

So how much work does that add to triangle setup? Setting up the $\frac{\lambda_0}{w}$ and $\frac{\lambda_1}{w}$ for the triangle requires 4 reciprocals, the triangle area (which we already computed for back-face culling!), and a few subtractions, multiplies and adds. Setting up the vertex attributes for interpolation is really cheap with the barycentric approach – two subtractions per attribute (if you don't use barycentric, you get some more multiply-add action here). Follow me? Probably not, unless you've implemented this before. Sorry about that – but it's fairly safe to ignore all this if you don't understand it.

Let's get back to why we're here: the one value we want to interpolate *right now* is Z, and because we computed Z as $\frac{z}{w}$ at the vertex level as part of projection (see previous part), so it's already divided by w and we can just interpolate it linearly in screen space. Nice. What we end up with is a plane equation for $Z = aX + bY + c$ that we can just plug X and Y into to get a value. So, here's the punchline of my furious hand-waving in the last few paragraphs: Interpolating Z at any given point boils down to two multiply-adds. (Starting to see why GPUs have fast multiply-accumulate units? This stuff is absolutely everywhere!).

Early Z/Stencil

Now, if you believe the place that graphics APIs traditionally put Z/Stencil processing into – right before alpha blend, way at the bottom of the pixel pipeline – you might be confused a bit. Why am I even discussing Z at the point in the pipeline where we are right now? We haven't even started shading pixels! The answer is simple: the Z and stencil tests reject pixels. Potentially the majority of them. You really, *really* don't want to completely shade a detailed mesh with complicated materials, to then throw away

95% of the work you just did because that mesh happens to be mostly hidden behind a wall. That's just a really stupid waste of bandwidth, processing power and energy. And in most cases, it's completely unnecessary: most shaders don't do anything that would influence the results of the Z test, or the values written back to the Z/stencil buffers.

So what GPUs actually do when they can is called "early Z" (as opposed to late Z, which is actually at the late stage in the pipeline that traditional API models generally display it at). This does exactly what it sounds like – execute the Z/stencil tests and writes early, right after the triangle has been rasterized, and before we start sending off pixels to the shaders. That way, we notice all the rejected pixels early, without wasting a lot of computation on them. However, we can't always do this: the pixel shader may ignore the interpolated depth value, and instead provide its own depth to be written to the Z-buffer (e.g. depth sprites); or it might use discard, alpha test, or alpha-to-coverage, all of which "kill" pixels/samples during pixel shader execution and mean that we can't update the Z-buffer or stencil buffer early because we might be updating depth values for samples that later get discarded in the shader!

So GPUs actually have two copies of the Z/stencil logic; one right after the rasterizer and in front of the pixel shader (which does early Z) and one after the shader (which does late Z). Note that we can still, in principle, do the depth testing in the early-Z stage even if the shader uses some of the sample-killing mechanism. It's only writes that we have to be careful with. The only case that really precludes us from doing any early Z-testing at all is when we write the output depth in the pixel shader – in that case the early Z unit simply has nothing to work with.

Traditionally, APIs just pretended none of this early-out logic existed; Z/Stencil was in a late stage in the original API model, and any optimizations such as early-Z had to be done in a way that was 100% functionally consistent with that model; i.e. drivers had to detect when early-Z was applicable, and could only turn it on when there were no observable differences. By now APIs have closed that gap; as of DX11, shaders can be declared as "force early-Z", which means they run with full early-Z processing even when the shader uses primitives that aren't necessarily "safe" for early-Z, and shaders that write depth can declare that the interpolated Z value is conservative (i.e. early Z reject can still happen).

Z/stencil writes: the full truth

Okay, wait. As I've described it, we now have two parts in the pipeline – early Z and late Z – that can both write to the Z/stencil buffers. For any given shader/render state combination that we look at, this will work – in the steady state. But that's not how it works in practice. What actually happens is that we render a few hundred to a few thousand batches per frame, switching shaders and render state regularly. Most of these shaders will allow early Z, but some won't. Switching from a shader that does early Z to one that does late Z is no problem. But going back from late Z to early Z is, if early Z does any writes: early Z is, well, earlier in the pipeline than late Z – that's the whole point! So we may start early-Z processing for one shader, merrily writing to the depth buffer while there's still stuff down in the pipeline for our old shader that's running late-Z and may be trying to write the same location at the same time – classic race condition. So how do we fix this? There's a bunch of options:

Once you go from early-Z to late-Z processing within a frame (or at least a sequence of operations for the same render target), you stay at late-Z until the next point where you flush the pipeline anyway. This works but potentially wastes lots of shader cycles while early-Z is unnecessarily off.

Trigger a (pixel) pipeline flush when going from a late-Z shader to an early-Z shader – also works, also not exactly subtle. This time, we don't waste shader cycles (or memory bandwidth) but stall instead – not much of an improvement.

But in practice, having Z-writes in two places is just bad news. Another option is to not ever write Z in the early-Z phase; always do the Z-writes in late-Z. Note that you need to be careful to make conservative Z-testing decisions during early Z if you do this! This avoids the race condition but means the early Z-test results may be stale because the Z-write for the currently-dispatched pixels won't happen until a while later.

Use a separate unit that keeps track of Z-writes for us and enforces the correct ordering; both early-Z and late-Z must go through this unit.

All of these methods work, and all have their own advantages and drawbacks. Again I'm not sure what current hardware does in these cases, but I have strong reason to believe that it's one of the last two options. In particular, we'll meet a functional unit later down the road (and the pipeline) that would be a good place to implement the last option.

But we're still doing all this testing per pixel. Can't we do better?

Hierarchical Z/Stencil

The idea here is that we can use our tile trick from rasterization again, and try to Z-reject whole tiles at a time, before we even descend down to the pixel level! What we do here is a strictly conservative test; it may tell us that “there might be pixels that pass the Z/stencil-test in this tile” when there are none, but it will never claim that all pixels are rejected when in fact they weren’t.

Assume here that we’re using “less”, “less-equal”, or “equal” as Z-compare mode. Then we need to store the *maximum* Z-value we’ve written for that tile, per tile. When rasterizing a triangle, we calculate the *minimum* Z-value the active triangle is going to write to the current tile (one easy conservative approximation is to take the min of the interpolated Z-values at the four corners of the current tile). If our triangle minimum-Z is larger than the stored maximum-Z for the current tile, the triangle is guaranteed to be completely occluded. That means we now need to track maximum-Z per-tile, and keep that value up to date as we write new pixels – though again, it’s fine if that information isn’t completely up to date; since our Z-test is of the “less” variety, values in the Z buffer will only get smaller over time. If we use a per-tile maximum-Z that’s a bit out of date, it just means we’ll get slightly worse early rejection rates than we could; it doesn’t cause any other problems.

The same thing works (with min/max and compare directions swapped) if we’re using one of the “greater”, “greater-equal” or “equal” Z-tests. What we can’t easily do is change from one of the “less”-based tests to a “greater”-based tests in the middle of the frame, because that would make the information we’ve been tracking useless (for less-based tests we need maximum-Z per tile, for greater-based tests we need minimum-Z per tile). We’d need to loop over the whole depth buffer to recompute min/max for all tiles, but what GPUs actually do is turn hierarchical-Z off once you do this (up until the next Clear). So: don’t do that.

Similar to the hierarchical-Z logic I’ve described, current GPUs also have hierarchical stencil processing. However, unlike hierarchical-Z, I haven’t seen much in the way of published literature on the subject (meaning, I haven’t run into it – there might be papers on it, but I’m not aware of them); as a game console developer you get access to low-level GPU docs which include a description of the underlying algorithms, but frankly, I’m definitely not comfortable writing about something here where really the only good sources I have are various GPU docs that came with a thick stack of NDAs. Instead I’ll just nebulously note that there’s magic pixie dust that can do certain kinds of stencil testing very efficiently under controlled circumstances, and leave you to ponder what that might be and how it might work, in the unlikely case that you deeply care about this – presumably because your father was killed by a hierarchical stencil unit and you’re now collecting information on its weak points for your revenge, or something like that.

Putting it all together

Okay, we now have all the algorithms and theory we need – let’s see how we can take our new set of toys and wire it up with what we already have!

First off, we now need to do some extra triangle setup for Z/attribute interpolation. Not much to be done about it – more work for triangle setup; that’s how it goes. After that’s coarse rasterization, which I’ve discussed in the previous part.

Then there’s hierarchical Z (I’m assuming less-style comparisons here). We want to run this between coarse and fine rasterization. First, we need the logic to compute the minimum Z estimates for each tile. We also need to store the per-tile maximum Zs, which don’t need to be exact: we can shave bits as long as we always round up! As usual, there’s a trade-off here between space used and early-rejection efficiency. In theory, you could put the Z-max info into regular memory. In practice, I don’t think anyone does this, because you want to make the hierarchical-Z decision without a ton of extra latency. The other option is to put dedicated memory for hierarchical Z onto the chip – usually as SRAM, the kind of memory you also make caches out of. For 24-bit Z, you probably need something like 10-14 bits per tile to store a reasonable-accuracy Z-max in a compact encoding. Assuming 8x8 tiles, that means less than 1MBit (128k) of SRAM to support resolutions up to 2048x2048 – sounds like a plausible order of magnitude to me. Note that these things are fixed size and shared for the whole chip; if you do a context switch, you lose. If you allocate the wrong depth buffers to this memory, you can’t use hierarchical Z on the depth buffers that actually matter, and you lose. That’s just how it goes. This kind of things is why hardware vendors regularly tell you to create your most important render targets and depth buffers first; they have a limited supply of this type of memory (there’s more like it, as you’ll see), and when it runs out, you’re out of luck. Note they don’t necessarily need to do this all-or-nothing; for example, if you have a really large depth buffer, you might only get hierarchical Z in the top left 2048x1536 pixels, because that’s how much fits into the Z-max memory. It’s not ideal, but still much better than disabling hierarchical-Z outright.

And by the way, “Real-Time Rendering” mentions at this point that “it is likely that GPUs are using hierarchical Z-buffers with more than two levels”. I doubt this is true, for the same reason that I doubt they use a multilevel hierarchical rasterizer: adding more levels makes the easy cases (large triangles) even faster while adding latency and useless work for small triangles: if you’re drawing a triangle that fits inside a single 8×8 tile, any coarser hierarchy level is pure overhead, because even at the 8×8 level, you’d just do one test to trivial-reject the triangle (or not). And again, for hardware, it’s not that big a performance issue; as long as you’re not consuming extra bandwidth or other scarce resources, doing more compute work than strictly necessary isn’t a big problem, as long as it’s within reasonable limits,

Hierarchical stencil is also there and should also happen prior to fine rast, most likely in parallel with hierarchical Z. We’ve established that this runs on air, love and magic pixie dust, so it doesn’t need any actual hardware and is probably always exactly right in its predictions. Ahem. Moving on.

After that is fine rasterization, followed in turn by early Z. And for early Z, there’s two more important points I need to make.

Revenge of the API order

For the past few parts, I’ve been playing fast and loose with the order that primitives are submitted in. So far, it didn’t matter; not for vertex shading, nor primitive assembly, triangle setup or rasterization. But Z is different. For Z-compare modes like “less” or “lessequal”, it’s very important what order the pixels arrive in; if we mess with that, we risk changing the results and introducing nondeterministic behavior. More importantly, as per the spec, we’re free to execute operations in any order *so long as it isn’t visible to the app*; well, as I just said, for Z processing, order is important, so we need to make sure that triangles arrive at Z processing in the right order (this goes for both early and late Z).

What we do in cases like this is go back in the pipeline and look for a reasonable spot to sort things into order again. In our current path, the best candidate location seems to be primitive assembly; so when we start assembling primitives from shaded vertex blocks, we make sure to assemble them strictly in the original order as submitted by the app to the API. This means we might stall a bit more (if the PA buffer holds an output vertex block, but it’s not the correct one, we need to wait and can’t start setting up primitives yet), but that’s the price of correctness.

Memory bandwidth and Z compression

The second big point is that Z/Stencil is a serious bandwidth hog. This has a couple of reasons. For one, this is the one thing we really run for all samples generated by the rasterizer (assuming Z/Stencil isn’t off, of course). Shaders, blending etc. all benefit from the early rejection we do; but even Z-rejected pixels do a Z-buffer read first (unless they were killed by hierarchical Z). That’s just how it works. The other big reason is that, when multisampling is enabled, the Z/stencil buffer is per sample; so 4x MSAA means 4x the memory bandwidth cost of Z? For something that takes a substantial amount of memory bandwidth even at no MSAA, that’s seriously bad news.

So what GPUs do is Z compression. There’s various approaches, but the general idea is always the same: assuming reasonably-sized triangles, we expect a lot of tiles to just contain one or maybe two triangles. If that happens, then instead of storing Z-values for the whole tile, we just store the plane equation of the triangle that filled up this tile. That plane equation is (hopefully) smaller than the actual Z data. Without MSAA, one tile covers 8×8 actual pixels, so triangles need to be relatively big to cover a full tile; but with 4x MSAA, a tile effectively shrinks to 4×4 pixels, and covering full tiles gets easier. There’s also extensions that can support 2 triangles etc., but for reasonably-sized tiles, you can’t go much larger than 2-3 tris and still actually save bandwidth: the extra plane equations and coverage masks aren’t free!

Anyway, point is: this compression, when it works, is fully lossless, but it’s not applicable to all tiles. So we need some extra space to denote whether a tile is compressed or not. We could store this in regular memory, but that would mean we now need to wait two full memory round-trips latencies to do a Z-read. That’s bad. So again, we add some dedicated SRAM that allows us to store a few (1-3) bits per tile. At its simplest, it’s just a “compressed” or “not compressed” flag, but you can get fancy and add multiple compression modes and such. A nice side effect of Z-compression is that it allows us to do fast Z-clears: e.g. when clearing to Z=1, we just set all tiles to “compressed” and store the plane equation for a constant Z=1 triangle.

All of the Z-compression thing, much like texture compression in the texture samplers, can be folded into memory access/caching logic, and made completely transparent to everyone else. If you don’t want to send the plane equations (or add the interpolator logic) to the Z memory access block, it can just infer them from the Z data and use some integer delta-coding scheme. This kind of

approach usually needs extra bits per sample to actually allow lossless reconstruction, but it can lead to simpler data paths and nicer interface between units, which hardware guys love.

And that's it for today! Next up: Pixel shading and what happens around it.

Postscript

As I said earlier, the topic of setting up interpolated attributes would actually make for a nice article on its own. I'm skipping that for now – might decide to fill this gap later, who knows.

Z processing has been in the 3D pipeline for ages, and a serious bandwidth issue for most of the time; people have thought long and hard about this problem, and there's a zillion tricks that go into doing "production-quality" Z-buffering for GPUs, some big, some small. Again, I'm just scratching the surface here; I tried to limit myself to the bits that are useful to know for a graphics programmer. That's why I don't spend much time on the details of hierarchical Z computations or Z compression and the like; all of this is very specific on hardware details that change slightly in every generation, and ultimately, mostly there's just no practical way you get to exploit any of this usefully: If a given Z-compression scheme works well for your scene, that's some memory bandwidth you can spend on other things. If not, what are you gonna do? Change your geometry and camera position so that Z-compression is more efficient? Not very likely. To a hardware designer, these are all algorithms to be improved on in every generation, but to a programmer, they're just facts of life to deal with.

This time, I'm not going into much detail on how memory accesses work in this stage of the pipeline. That's intentional. There's a key to high-throughput pixel shading and other per-pixel or per-sample processing, but it's later in the pipeline, and we're not there yet. Everything will be revealed in due time :)

From → Coding, Graphics Pipeline

13 Comments

1. KeyJ permalink

Why is the order of the primitives so important and needs to be preserved? As long as only fully opaque triangles are rendered (i.e. blending is off, which should be the case for the largest part of typical scenes), the rendering order shouldn't make any difference (except maybe some additional overdraw if the application sorted primitives by depth), or am I mistaken here?

Reply

- o [fgiesen permalink](#)

Suppose you're rendering a batch that has 3 triangles, all with the exact same vertex coordinates but different colors (this example is a bit contrived, but nevertheless completely valid). Let's say the first triangle is red, the second green, and the third blue.

If you render this with "less" as compare mode onto a freshly cleared render target and depth buffer, you will only see the red triangle. If you use "less_equal", you will only see the blue triangle. That's what the functional API model (that is fully sequential) requires.

Allowing reordering of triangles in the pipeline, even just for opaque rendering, breaks this: for both compare modes, you can make any of the 3 triangles come out "on top" using an appropriate ordering.

This is a bit tricky, since you *can* allow triangle reordering within a batch if you only care about the contents of the Z-buffer (e.g. for shadow map rendering with a NULL pixel shader). For both compare functions, the Z-buffer at the end of the frame will contain at each location the minimum of all depth values that have been written there.

However, the output of our Z-processing isn't just the Z-buffer, it's also an updated coverage mask that tells us which pixels/samples to shade. And for the coverage mask computation, order matters.

Reply

2. TomF permalink

I believe it's common in HW to have all real Z done late. That result is then fed back "upstream" into the coarse Z unit, and the only sort of early Z is coarse Z. This is conservative and can be "late" – a triangle occluded by the previous triangle will still be shaded – but it is always safe and requires no mode-switching shenanigans.

Reply

- o [fgiesen permalink](#)

I'm not sure which type of implementation is more common these days, but I've definitely worked with chips that support both pre- and post-shading Z (including ones that don't do any coarse Z). But I've never seen anyone use anything but the obvious brute-force solution for the late Z->early Z transition: just flush the pixel pipeline when you're doing that switch.

Reply

3. piyush permalink

contrary to popular belief, computer scholars do have a sense of humor.
and you, sir are a living example

Reply

4. Martin Wardener permalink

"early Z is, well, earlier in the pipeline than early Z" .. I assume you mean "..than late Z" ..?

Reply

- o [fgiesen permalink](#)

Indeed! Thank you, fixed.

Reply

5. PixInverse permalink

I'm little confused by your explanation of perspective correct interpolation. Once we have performed the perspective divide $x/w, y/w, z/w$, then $Z(z/w)$, the value we want to interpolate, doesn't vary linearly anymore across the surface of the 2D triangle. What now varies linearly is $1/Z$.

To interpolate a vertex attribute correctly, we first need to divide by the vertex attribute value by Z of the vertex it is defined to, then linearly interpolate them, and then finally multiply the result by Z -Interpolated, which is the depth of the point on the triangle, that the pixel overlaps.

$$Z_{\text{correct}} = 1 / (\text{Interpolated } 1/Z)$$

$$U_{\text{correct}} = (\text{Interpolated } u/Z) * Z_{\text{correct}}$$

$$V_{\text{correct}} = (\text{Interpolated } v/Z) * Z_{\text{correct}}$$

In your post above, you says $u/w, v/w$ and z/w can be interpolated linearly for perspective correct interpolation. Could you please explain if I'm missing anything here?

Reply

- o [fgiesen permalink](#)

$Z=z/w$ ends up being an affine function of screen-space X and Y and is what's used for depth buffering.

It's never used for interpolation. Remember Z is set up by the projection matrix to reach (depending on the convention) either 0 and 1 or -1 and 1 at z_{near} and z_{far} , respectively.

What you use for interpolation is $1/w$ (which is also linear in screen space). Normally you set up perspective-corrected barycentric coordinates (usually called I and J) and then you interpolate $I/w, J/w$ and $1/w$ in screen space, and then for every pixel you solve $I=(I/w)/(1/w), J=(J/w)/(1/w)$ and use I and J to interpolate attributes. (As many as you want, and without having to individually divide them all through by the vertex w's).

Reply

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog
2. A trip through the Graphics Pipeline 2011, part 8 « The ryg blog
3. A trip through the Graphics Pipeline 2011, part 9 « The ryg blog
4. Linear Depth | The Devil In The Details

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 8

July 10, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

In this part, I'll be dealing with the first half of pixel processing: dispatch and actual pixel shading. In fact, this is really what most graphics programmer think about when talking about pixel processing; the alpha blend and late Z stages we'll encounter in the next part seem like little more than an afterthought. In hardware, the story is a bit more complicated, as we'll see – there's a reason I'm splitting pixel processing into two parts. But I'm getting ahead of myself. At the point where we're entering this stage, the coordinates of pixels (or, actually, quads) to shade, plus associated coverage masks, arrive from the rasterizer/early-Z unit – with triangle in the exact same order as submitted by the application, as I pointed out last time. What we need to do here is to take that linear, sequential stream of work and farm it out to hundreds of shader units, then once the results are back, we need to merge it back into one linear stream of memory updates.

That's a textbook example of fork/join-parallelism. This part deals with the fork phase, where we go wide; the next part will explain the join phase, where we merge the hundreds of streams back into one. But first, I have a few more words to say about rasterization, because what I just told you about there being just one stream of quads coming in isn't quite true.

Going wide during rasterization

To my defense, what I told you *used* to be true for quite a long time, but it's a serial part of the pipeline, and once you throw in excess of 300 shader units at a problem, serial parts of the pipeline have the tendency to become bottlenecks. So GPU architects started using multiple rasterizers; as of 2010, **NVidia employs four rasterizers**

(http://www.highperformancegraphics.org/previous/www_2010/media/Hot3D/HPG2010_Hot3D_NVIDIA.pdf) and **AMD uses two** (http://www.highperformancegraphics.org/previous/www_2010/media/Hot3D/HPG2010_Hot3D_AMD.pdf). As a side note, the NV presentation also has a few notes on the requirement to keep stuff in API order. In particular, you really do need to sort primitives back into order prior to rasterization/early-Z, like I mentioned last time; doing it just before alpha blend (as you might be inclined to do) doesn't work.

The work distribution between rasterizers is based on the tiles we've already seen for early-Z and coarse rasterization. The frame buffer is divided into tile-sized regions, and each region is assigned to one of the rasterizers. After setup, the bounding box of the triangle is consulted to find out which triangles to hand over to which rasterizers; large triangles will always be sent to all rasterizers, but smaller ones can hit as little as one tile and will only be sent to the rasterizer that owns it.

The beauty of this scheme is that it only requires changes to the work distribution and the coarse rasterizers (which traverse tiles); everything that only sees individual tiles or quads (that is, the pipeline from hierarchical Z down) doesn't need to be modified. The problem is that you're now dividing jobs based on screen locations; this can lead to a severe load imbalance between the rasterizers (think a few hundred tiny triangles all inside a single tile) that you can't really do anything about. But the nice thing is that everything that adds ordering constraints to the pipeline (Z-test/write order, blend order) comes attached to specific frame-buffer locations, so screen-space subdivision works without breaking API order – if this wasn't the case, tiled renderers wouldn't work.

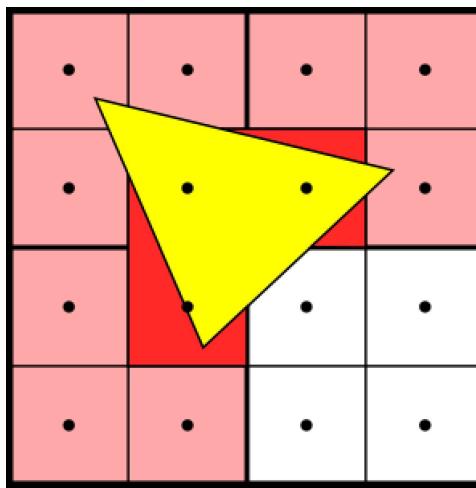
You need to go wider!

Okay, so we don't get just one linear stream of quad coordinates plus coverage masks in, but between two and four. We still need to farm them out to hundreds of shader units. It's time for another dispatch unit! Which first means another buffer. But how big are the batches we send off to the shaders? Here I go with NVidia figures again, simply because they mention this number in **public white papers**

(http://www.nvidia.com/content/PDF/fermi_white_papers/NVIDIA_Fermi_Compute_Architecture_Whitepaper.pdf); AMD probably also states that information somewhere, but I'm not familiar with their terminology for it so I couldn't do a direct search

for it. Anyway, for NVidia, the unit of dispatch to shader units is 32 threads, which they call a “Warp”. Each quad has 4 pixels (each of which in turn can be handled as one thread), so for each shading batch we issue, we need to grab 8 incoming quads from the rasterizer before we can send off a batch to the shader units (we might send less in case there’s a shader switch or pipeline flush).

Also, this is a good point to explain why we’re dealing with quads of 2×2 pixels and not individual pixels. The big reason is derivatives. Texture samplers depend on screen-space derivatives of texture coordinates to do their mip-map selection and filtering (as we saw back in [part 4](https://fgiesen.wordpress.com/2011/07/04/a-trip-through-the-graphics-pipeline-2011-part-4/) (<https://fgiesen.wordpress.com/2011/07/04/a-trip-through-the-graphics-pipeline-2011-part-4/>)); and, as of shader model 3.0 and later, the same machinery is directly available to pixel shaders in the form of derivative instructions. In a quad, each pixel has both a horizontal and vertical neighbor within the same quad; this can be used to estimate the derivatives of parameters in the x and y directions using [finite differencing](http://en.wikipedia.org/wiki/Finite_difference) (http://en.wikipedia.org/wiki/Finite_difference) (it boils down to a few subtractions). This gives you a very cheap way to get derivatives at the cost of always having to shade groups of 2×2 pixels at once. This is no problem in the interior of large triangles, but means that between 25-75% of the shading work for quads generated for triangle edges is wasted. That’s because all pixels in a quad, even the masked ones, get shaded. This is necessary to produce correct derivatives for the pixels in the quad that are visible. The invisible but still-shaded pixels are called “helper pixels”. Here’s an illustration for a small triangle:



(https://fgiesen.files.wordpress.com/2011/07/quad_coverage.png)

The triangle intersects 4 quads, but only generates visible pixels in 3 of them. Furthermore, in each of the 3 quads, only one pixel is actually covered (the sampling points for each pixel region are depicted as black circles) – the pixels that are filled are depicted in red. The remaining pixels in each partially-covered quad are helper pixels, and drawn with a lighter color. This illustration should make it clear that for small triangles, a large fraction of the total number of pixels shaded are helper pixels, which has attracted some [research attention](http://graphics.stanford.edu/papers/fragmerging/shade_sig10.pdf) (http://graphics.stanford.edu/papers/fragmerging/shade_sig10.pdf) on how to merge quads of adjacent triangles. However, while clever, such optimizations are not permissible by current API rules, and current hardware doesn’t do them. Of course, if the HW vendors at some point decide that wasted shading work on quads is a significant enough problem to force the issue, this will likely change.

Attribute interpolation

Another unique feature of pixel shaders is attribute interpolation – all other shader types, both the ones we’ve seen so far (VS) and the ones we’re still to talk about (GS, HS, DS, CS) get inputs directly from a preceding shader stage or memory, but pixel shaders have an additional interpolation step in front of them. I’ve already talked a bit about this in the [previous part](https://fgiesen.wordpress.com/2011/07/08/a-trip-through-the-graphics-pipeline-2011-part-7/) (<https://fgiesen.wordpress.com/2011/07/08/a-trip-through-the-graphics-pipeline-2011-part-7/>) when discussing Z, which was the first interpolated attribute we saw.

Other interpolated attributes work much the same way; a plane equation for them is computed during triangle setup (GPUs may choose to defer this computation somewhat, e.g. until it’s known that at least one tile of the triangle passed the hierarchical Z-test, but that shall not concern us here), and then during pixel shading, there’s a separate unit that performs attribute interpolation using the pixel positions of the quads and the plane equations we just computed.

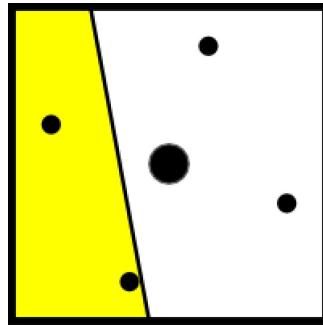
Update: Marco Salvi points out (in the comments below) that while there used to be dedicated interpolators, by now the trend is towards just having them return the barycentric coordinates to plug into the plane equations. The actual evaluation (two multiply-adds per attribute) can be done in the shader unit.

All of this shouldn't be surprising, but there's a few extra interpolation types to discuss. First, there's "constant" interpolators, which are (surprise!) constant across the primitive and take the value for each vertex attribute from the "leading vertex" (which vertex that is is determined during primitive setup). Hardware may either have a fast-path for this or just set up a corresponding plane equation; either way works fine.

Then there's no-perspective interpolation. This will usually set up the plane equations differently; the plane equations for perspective-correct interpolation are set up either for X, Y-based interpolation by dividing the attribute values at each vertex by the corresponding w, or for barycentric interpolation by building the triangle edge vectors. Non-perspective interpolated attributes, however, are cheapest to evaluate when their plane equation is set up for X, Y-based interpolation *without* dividing the values at each vertex by the corresponding w.

"Centroid" interpolation is tricky

Next, we have "centroid" interpolation. This is a flag, not a separate mode; it can be combined both with the perspective and no-perspective modes (but not with constant interpolation, because it would be pointless). It's also terribly named and a no-op unless multisampling is enabled. With multisampling on, it's a somewhat hacky solution to a real problem. The issue is that with multisampling, we're evaluating triangle *coverage* at multiple sample points in the rasterizer, but we're only doing the actual *shading* once per pixel. Attributes such as texture coordinates will be interpolated at the pixel center position, as if the whole pixel was covered by the primitive. This can lead to problems in situations such as this:



(https://fgiesen.files.wordpress.com/2011/07/msaa_samples.png)

Here, we have a pixel that's partially covered by a primitive; the four small circles depict the 4 sampling points (this is the default 4x MSAA pattern) while the big circle in the middle depicts the pixel center. Note that the big circle is outside the primitive, and any "interpolated" value for it will actually be linear extrapolation; this is a problem if the app uses texture atlases, for example. Depending on the triangle size, the value at the pixel center can be very far off indeed. Centroid sampling solves this problem. The original explanation was that the GPU takes all of the samples covered by the primitive, computes their centroid, and samples at that position (hence the name). This is usually followed by the addition that this is just a conceptual model, and GPUs are free to do it differently, so long as the point they pick for sampling is within the primitive.

If you think it somewhat unlikely that the hardware actually counts the covered samples, sums them up, then divides by the count, then join the club. Here's what *actually* happens:

If all sample points cover the primitive, interpolation is done as usual, i.e. at the pixel center (which happens to be the centroid of all sample positions for all reasonable sampling patterns).

If not all sample points cover the triangle, the hardware picks one of the sample points that do and evaluates there. All covered sample points are (by definition) inside the primitive so this works.

That picking used to be arbitrary (i.e. left to the hardware); I believe by now DX11 actually prescribes exactly how it's done, but this more a matter of getting consistent results between different pieces of hardware than it is something that API users will actually care about. As said, it's a bit hacky. It also tends to mess up derivative calculations for quads that have partially covered pixels – tough luck. What can I say, it may be industrial-strength duct tape, but it's still duct tape.

Finally (new in DX11!) there's "pull-model" attribute interpolation. Regular attribute interpolation is done automatically before the pixel shader starts; pull-model interpolation adds actual instructions that do the interpolation to the pixel shader. This allows the shader to compute its own position to sample values at, or to only interpolate attributes in some branches but not in others. What it boils down to is the pixel shader being able to send additional requests to the interpolation unit while the shader is running.

The actual shader body

Again, the general shader principles are well-explained in the API documentation, so I'm not going to talk about how individual instructions work; generally, the answer is "as you would expect them to". There are however some interesting bits about pixel shader execution that are worth talking about.

The first one is: texture sampling! Wait, didn't I wax on about texture samplers for quite some time in part 4 already? Yes, but that was the texture sampler side of things – and if you remember, there was that one bit about texture cache misses being so frequent that samplers are usually designed to sustain at least one miss to main memory per request (which is 16-32 pixels, remember!) without stalling. That's a lot of cycles – hundreds of them. And it would be a tremendous waste of perfectly good ALUs to keep them idle while all this is going on.

So what shader units actually do is switch to a different batch after they've issued a texture sample; then when that batch issues a texture sample (or completes), it switches back to one of the previous batches and checks if the texture samples are there yet. As long as each shader unit has a few batches it can work on at any given time, this makes good use of available resources. It does increase latency for completion of individual batches though – again, a latency-vs-throughput trade-off. By now you should know which side wins on GPUs: Throughput! *Always*. One thing to note here is that keeping multiple batches (or "Warps" on NVidia hardware, or "Wavefronts" for AMD) running at the same time requires more registers. If a shader needs a lot of registers, a shader unit can keep less warps around; and if there are less of them, the chance that at some point you'll run out of runnable batches that aren't waiting on texture results is higher. If there's no runnable batches, you're out of luck and have to stall until one of them gets its results back. That's unfortunate, but there's limited hardware resources for this kind of thing – if you're out of memory, you're out of memory, period.

Another point I haven't talked about yet: Dynamic branches in shaders (i.e. loops and conditionals). In shader units, work on all elements of each batch usually proceeds in lockstep. All "threads" run the same code, at the same time. That means that ifs are a bit tricky: If *any* of the threads want to execute the "then"-branch of an if, all of them have to – even though most of them may end up ignoring the results using a technique called **predication** (http://en.wikipedia.org/wiki/Branch_predication), because they didn't want to descend down there in the first place.. Similarly for the "else" branch. This works great if conditionals tend to be coherent across elements, and not so great if they're more or less random. Worst case, you'll always execute both branches of every if. Ouch. Loops work similarly – as long as at least one thread wants to keep running a loop, all of the threads in that batch/Warp/Wavefront will.

Another pixel shader specific is the **discard** instruction. A pixel shader can decide to "kill" the current pixel, which means it won't get written. Again, if all pixels inside a batch get discarded, the shader unit can stop and go to another batch; but if there's at least one thread left standing, the rest will be dragged along. DX11 adds more fine-grained control here by way of writing the output pixel coverage from the pixel shader (this is always ANDed with the original triangle/Z-test coverage, to make sure that a shader can't write outside its primitive, for sanity). This allows the shader to discard individual samples instead of whole pixels; it can be used to implement Alpha-to-Coverage with a custom dithering algorithm in the shader, for example.

Pixel shaders can also write the output depth (this feature has been around for quite some time now). In my experience, this is an excellent way to shoot down early-Z, hierarchical Z and Z compression and in general get the slowest path possible. By now, you know enough about how these things work to see why. :)

Pixel shaders produce several outputs – in general, one 4-component vector for each render target, of which there can be (currently) up to 8. The shader then sends the results on down the pipeline towards what D3D calls the "Output Merger". This'll be our topic next time.

But before I sign off, there's one final thing that pixel shaders can do starting with D3D11: they can write to Unordered Access Views (UAVs) – something which only compute and pixel shaders can do. Generally speaking, UAVs take the place of render targets during compute shader execution; but unlike render targets, the shader can determine the position to write to itself, and there's no implicit API order guarantee (hence the "unordered access" part of the name). For now, I'll only mention that this functionality exists – I'll talk more about it when I get to Compute Shaders.

Update: In the comments, Steve gave me a heads-up about the correct AMD terminology (the first version of the post didn't have the "Wavefronts" name because I couldn't remember it) and also posted a link to [this great presentation by Kayvon Fatahalian](http://bps10.idav.ucdavis.edu/talks/03-fatahalian_gpuArchTeraflop_BPS_SIGGRAPH2010.pdf) (http://bps10.idav.ucdavis.edu/talks/03-fatahalian_gpuArchTeraflop_BPS_SIGGRAPH2010.pdf) that explains shader execution on GPUs, with a lot more pretty pictures that I can be bothered to make :). You should really check it out if you're interested in how shader cores work.

And... that's it! No big list of caveats this time. If there's something missing here, it's because I've genuinely forgotten about it, not because I decided it was too arcane or specific to write up here. Feel free to point out omissions in the comments and I'll see what I can do.

From → Coding, Graphics Pipeline

19 Comments

1. Marco Salvi permalink

I am really enjoying your series. You should collect it in a single document once it's complete.

A minor note on attributes interpolation:

".. then during pixel shading, there's a separate unit that performs attribute interpolation using the pixel positions of the quads and the plane equations we just computed."

This used to be true. These days you are mostly likely to find a separate unit that "simply" generates barycentric coordinates (according the type of interpolation requested), while the rest of the per-attribute calculation is performed on the shader cores.

Reply

- [fgiesen permalink](#)

"I am really enjoying your series. You should collect it in a single document once it's complete."

Thanks! Once I'm done I'll at least go over the text again, revise it a bit and probably add a few more illustrations. I may collect it as a single document – not decided yet.

"These days you are mostly likely to find a separate unit that "simply" generates barycentric coordinates (according the type of interpolation requested), while the rest of the per-attribute calculation is performed on the shader cores."

Ah. I know that the full-blown interpolators used to be there, but it's generally not something that gets much attention in official presentations, so it's hard to pinpoint when exactly they switch :)

Makes sense, though. I suspect that everything that involves fused multiply-adds will end up getting moved into the shader units, sooner or later. Especially as compute to texture ratio keeps going up steadily, there's no reason not to.

Reply

2. Rex Guo permalink

Thanks for another great write-up, ryg!

While on the topic of pixel/fragment shaders, could you give some guidance on the ideal ratio of compute/memory instruction mix? Is it still ~20:1 ? Does NV also publish the cycle count for the math ops?

Reply

- [fgiesen permalink](#)

"could you give some guidance on the ideal ratio of compute/memory instruction mix"

That kind of information is outside the scope of this series; it depends too much on the underlying hardware. That said, as you get more transistors with lower power consumption, it's easier to add compute power than it is to add (DRAM) memory bandwidth. And memory latency is worse – lowering latency without harming bandwidth or power consumption is seriously hard.

NV does publish the latency of their math ops – it's in the same document I got the memory access latency from (on the same page, too!). They use fully pipelined ALUs with very high latencies; 18+ cycles per operation. This means that they do need to switch warps every cycle – you need 18 warps running concurrently on a shader unit to not ever hit instruction dependency stalls, though if the shaders have several independent dependency chains, that number goes down. If you compare the ~20 cycles for ALU ops to the ~400 cycles for memory access, you can see where the 20:1 ratio comes from – while a warp is waiting on a texture fetch, it temporarily drops out of the rotation, so you need correspondingly more runnable warps to fill the gap.

Reply

3. Steve permalink

Excellent posts, thanks!

> "Warps", or whatever the AMD terminology is

AFAIK AMD calls them "wavefronts" ("Running Code at a Teraflop: How a GPU Shader Core Works", <http://bps10.idav.ucdavis.edu/> – a great companion course for this stuff)

Reply

- o [fgiesen permalink](#)

Ah, Kayvons presentation is one I've been looking for, but I couldn't remember the title or the name of the author, which makes it hard to find :). I'll update the article accordingly, thanks a bunch!

Reply

4. Barbie permalink

Do you have any insight about the derivative computation based on helper pixels (outside of the triangle) ? I've never really looked into the artifacts it can trigger, but I assume e.g. computing texture coordinates based on a previous texture lookup can end up looking really weird if the texture lookup ended up wrapping around.

Reply

- o [fgiesen permalink](#)

Extrapolation for helper pixels doesn't cause problems with linear attributes: if you think of it in world space, the texture coordinates are interpolated linearly across the triangle, so their world-space derivatives are constants. In screen-space (post-projection), there's an added perspective distortion, which is a projective transform that takes the plane of the triangle to the plane the screen quad lies in (i.e. viewing plane). Since it's a transform between planes, it's well-defined even outside the triangle, and approximating the derivatives with screen-space finite differences isn't any more (or less) correct outside the tri than it is inside of it.

If your quantities are not linear, or cross a seam, then yes, that causes visible artifacts. Usually you get some 2x2 quads that are a lot blurrier than others (because some quads see artificially high derivatives, which makes them pick small mip levels). Shader authors by now are usually aware of the problem, but if you look for this in games, you'll find it :). Early- to mid-2000s shaders for things like Environment-Mapped Bump Mapping have an especially bad case of this (because they were written for 1.x or 2.x shader HW which doesn't give you any means of specifying gradients directly – even if you had enough instruction slots to compute them...)

Reply

5. Egor Yusov permalink

"In a quad, each pixel has both a horizontal and vertical neighbor within the same quad; this can be used to estimate the derivatives of parameters in the x and y directions using finite differencing (it boils down to a few subtractions)"

Does this mean that for instance both pixels in each row get the same horizontal derivative?

It seems like computing just one difference should not be a problem for GPU, while different sources tell us that computing derivatives is quite expensive and suggest not computing them in a shader, when possible. There are even special instructions that compute coarse derivatives. Do you know how these instructions work?

Reply

- o [fgiesen permalink](#)

Yes. The x deltas and y deltas are the same for the 2 pixels in a row or column, respectively. And this is with "fine" derivatives – with "coarse" derivatives, both x and y derivatives are the same for all pixels in a quad.

Microsoft have updated the online D3D11 docs to include the instruction specs, so here they are for reference:

Fine derivatives

Coarse derivatives

The computation itself is cheap, but the necessary shuffling to get the difference between two values computed in different threads usually has some extra cost. While this might make it more expensive than other ALU instructions, it's not particularly bad, and certainly much cheaper than e.g. texture fetches.

While the name might suggest otherwise, there should be no appreciable speed difference between computing coarse and fine derivatives.

Reply

- o [fgiesen permalink](#)

For what it's worth, on implementation of derivative instructions in hardware: By now there's specs for both recent AMD and NV ISAs (or at least a reasonably close approximation thereof) online, so you can just look it up.

In NVidias PTX 3.1 ISA, you could implement a coarse x derivative like this: (hope I got this right – this is the right sequence of operations, but I might have screwed up the masks)

```
shfl.idx.b32 Rtl, Rsrc, 0x00, 0x1c00;
shfl.idx.b32 Rtr, Rsrc, 0x01, 0x1c00;
sub.f32      Rdsrc_dx, Rtr, Rtl;
```

Coarse y:

```
shfl.idx.b32 Rtl, Rsrc, 0x00, 0x1c00;
shfl.idx.b32 Rbl, Rsrc, 0x02, 0x1c00;
sub.f32      Rdsrc_dy, Rbl, Rtl;
```

Fine x:

```
shfl.bfly.b32 Ropposite|p, Rsrc, 0x01, 0x1e00;
@p sub.f32     Rdsrc_dx, Ropposite, Rsrc;
@!p sub.f32    Rdsrc_dx, Rsrc, Ropposite;
```

Fine y:

```
shfl.bfly.b32 Ropposite|p, Rsrc, 0x02, 0x1d00;
@p sub.f32     Rdsrc_dy, Ropposite, Rsrc;
@!p sub.f32    Rdsrc_dy, Rsrc, Ropposite;
```

With AMD's Southern Islands ISA (used for their recent GCN cores), you can perform the shuffling using the `DS_SWIZZLE_B32` instruction and the subtractions are simple. On both architectures works out to roughly three times the cost of a regular ALU operation (add, sub, mad, etc.). Not free, but not particularly expensive either.

6. Niad permalink

Hi thank you for creating this series!

I was wondering about the quad overshading.

If the derivatives are only required for texturing and calls to `ddx/ddy`, why not disable it if a shader does not use these features?

Are derivatives required for anything else? Does current hardware disable the quads if they aren't used?

This might seem random, but I have a project that works exactly like this, no texturing, very small triangles. It seems be bottlenecked by this issue.

Reply

- o [fgiesen](#) permalink

All current 3D hardware that I know of uses quads always. It is not something you can just disable.

This is not a software/driver thing; quad granularity is part of the design of most fixed-function blocks in the pipeline that are used at the per-pixel level. the rasterizer determines coverage in terms of quads, not pixels. Depth and stencil testing is done at quad (or larger) granularity at a time. Pixel/fragment shading is done on groups of pixels at a time (usually between 16 and 64), and the hardware that determines these groups works on quads, not individual pixels. Attributes are interpolated based on the triangle that a fragment came from. With quad-based shading and 64-wide wavefronts, that means that a wavefront can reference up to 16 unique triangles; without quad-based shading, it would be 64. There needs to be storage for the plane equations of attributes, and hardware to set up these plane equations for interpolation per triangle from vertex attributes before the shader runs. Getting rid of quads would mean you would need 4x the amount of storage for plane equations (to handle the worst case), 4x more (or 4x faster) hardware to set up these attributes, and so forth. And once shading is done, blending and writing to memory is usually done on a quad basis as well.

It's not that it's impossible to design HW that doesn't use quads; but supporting true pixel-granularity shading essentially boils down to re-designing a bunch of hardware blocks to have one or more of: 4x the clock rate, 4x the operation width, 4x the amount of routing/control logic, 4x the amount of buffer space. It's not something you just "enable" – it's a major

change, costly, power-hungry, it doesn't benefit workloads that GPU vendors (currently) care about, and is thus unlikely to get implemented any time soon.

Reply

7. [nlguillemot permalink](#)

minor typo: "With multisampling ob" should be "With multisampling on", I assume.

Reply

- o [fgiesen permalink](#)

Thanks, fixed!

Reply

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog
2. Viaje alucinante por un pipeline grafico « martin b.r.
3. A trip through the Graphics Pipeline 2011, part 13 « The ryg blog
4. Optimizing the basic rasterizer « The ryg blog

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 9

July 12, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Welcome back! This post deals with the second half of pixel processing, the “join phase”. The previous phase was all about taking a small number of input streams and turning them into lots of independent tasks for the shader units. Now we need to fold that large number of independent computations back into one (correctly ordered) stream of memory operations. As I already did in the posts on rasterization and early Z, I’ll first give a quick description of what needs to be done on a general level, and then I’ll go into how this is mapped to hardware.

Merging pixels again: blend and late Z

At the bottom of the pipeline (in what D3D calls the “Output Merger” stage), we have late Z/stencil processing and blending. These two operations are both relatively simple computationally, and they both update the render target(s) / depth buffer respectively. “Update” operation here means they’re of the read-modify-write variety. Because all of this happens for every quad that makes it this far through the pipeline, it’s also bandwidth-intensive. Finally, it’s order-sensitive (both blending and Z processing need to happen in API order), so we need to make sure to sort processed quads into order first.

I’ve already explained Z-processing, and blending is one of these things that work pretty much as you’d expect; it’s a fixed-function block that performs a multiply, a multiply-add and maybe some subtractions first, per render target. This block is kept deliberately simple; it’s separate from the shader units so it needs its own ALU, and we’d really prefer for it to be as small as possible: we want to spend our chip area (and power budget) on ALUs in the shader units, where they benefit every code that runs on the GPU, not on a fixed-function unit that’s only used at the end of the pixel pipeline. Also, we need it to have a short, predictable latency: this part of the pipeline needs to process data in-order to be correct. This limits our options as far as trading throughput for latency is concerned; we can still process quads that don’t overlap in parallel, but if we e.g. draw lots of small triangles, we’ll have multiple quads coming in for every screen location, and we’d better be able to write them out as quickly as they come, or else all our massively parallel pixel processing was for nought.

Meet the ROPs

ROPs are the hardware units that handle this part of the pipeline (as you can tell by the plural, there’s more than one). The acronym, depending on who you ask, stands for “Render OutPut unit”, “Raster Operations Pipeline”, or “Raster Operations Processor”. The actual name is fairly archaic – it derives from the days of pure 2D hardware acceleration, with hardware whose main purpose was to do fast **Bit blits** (http://en.wikipedia.org/wiki/Bit_blt). The classic 2D ROP design has three inputs – the current (destination) pixel value in the frame buffer, the source data, and a mask input – then computes some function of the 3 values and writes the results back to the frame buffer. Note this is before true color displays: the image data was usually in bit plane format and the function was some binary logic function. Then at some point bit planes died out (in favor of “chunky” representations that keep the bits for a pixel together), true color became the norm, the on-off mask was replaced with an alpha channel and the bitwise operations with blends, but the name stuck. So even now in 2011, when about the last remnant of that original architecture is the “logic op” in OpenGL, we still call them ROPs.

So what do we need to do, in hardware, for blend/late Z? A simple plan:

1. Read original render target/depth buffer contents from memory – memory access, long latency. Might also involve depth buffer and render target decompression! (I’ll explain render target compression later)
2. Sort incoming shaded quads into the right (API) order. This takes some buffering so we don’t immediately stall when quads don’t finish in the right order (think loops/branches, discard, and variable texture fetch latency). Note we only need to sort

- based on primitive ID here – two quads from the same primitive can never overlap, and if they don't overlap they don't need to be sorted!
3. Perform the actual blend/late Z/stencil operation. This is math – maybe a few dozen cycles worth, even with deeply pipelined units.
 4. Write the results back to memory again, compressing etc. along the way – long latency again, though this time we're not waiting for results so it's less of a problem at this end.

So, build the late-Z/blending unit, add some compression logic, wire it up to memory on one side and do some buffering of shaded quads on the other side and we're done, right?

Well, in theory anyway.

Except we need to cover the long latencies somehow. And all this happens for *every single pixel* (well, quad, actually). So we need to worry about memory bandwidth too... memory bandwidth? Wasn't there something about memory bandwidth? Watch closely now as I pull a bunny out of a hat after I put it there way back in **part 2** (<https://fgiesen.wordpress.com/2011/07/02/a-trip-through-the-graphics-pipeline-2011-part-2/>) (uh oh, that was more than a week ago – hope that critter is still OK in there...).

Memory bandwidth redux: DRAM pages

In part 2, I described the 2D layout of DRAM, and how it's faster to stay within a single row because changing the active row takes time – so for ideal bandwidth you want to stay in the same row between accesses. Well, the thing is, single DRAM rows are kinda large. Individual DRAM chips go up into the Gigabit range in size these days, and while they're not necessarily square (in fact a 2:1 aspect ratio seems to be preferred), you can still do a rough calculation of how many rows and columns there would be; for 512 Megabit (=64MB), we'd expect something like 16384×32768, i.e. a single row is about 32k bits or 4k bytes (or maybe 2k, or 8k, but somewhere in that ballpark – you get the idea). That's a rather inconvenient size to be making memory transactions in.

Hence, a compromise: the page. A DRAM page is some more conveniently sized slice of a row (by now, usually 256 or 512 bits) that's commonly transferred in a single burst. Let's take 512 bits (64 bytes) for now. At 32 bits per pixel – the standard for depth buffers and still fairly common for render targets although rendering workloads are definitely shifting towards 64 bit/pixel formats – that's enough memory to fit data for 16 pixels in. Hey, that's funny – we're usually shading pixels in groups of 16 to 64! (NV is a bit closer to the smaller end, AMD favors the larger counts). In fact, the 8x8 tile size I've been quoting in the rasterizer / early Z parts comes from AMD; I wouldn't be surprised if NV did coarse traversal (and hierarchical Z, which they dub "Z-cull") on 4x4 tiles, though a quick web search turned up nothing to either confirm this or rule it out. Either way, the plot thickens. Could it be that we're trying to traverse pixels in an order that gives good DRAM page coherency? You bet we are. Note that this has implications for internal render target layout too: we want to make sure pixels are stored such that a single DRAM page actually has a useful shape; for shading purposes, a 4x4 or 8x2 pixel DRAM page is a lot more useful than a 16x1 pixel one (remember – quads). Which is why render targets usually don't have a fully linear layout in memory.

That gives us yet another reason to shade pixels in groups, and also yet another reason to do a two-level traversal. But can we milk this some more? You bet we can: we still have the memory latency to cover. Usual disclaimer: This is one of the places where I don't have detailed information on what GPUs actually do, so what I'm describing here is a guess, not a fact. Anyway, as soon as we've rasterized a tile, we know whether it generates any pixels or not. At that point, we can select a ROP to handle our quads for that tile, and queue a command to fetch the associated frame buffer data into a buffer. By the point we get shaded quads back from the shader units, that data should be there, and we can start blending without delay (of course, if blending is off or identity, we can skip this load altogether). Similarly for Z data – if we run early Z before the pixel shader, we might need to allocate a ROP and fetch depth/stencil data earlier, maybe as soon as a tile has passed the coarse Z test. If we run late Z, we can just prefetch the depth buffer data at the same time we grab the framebuffer pixels (unless Z is off completely, that is).

All of this is early enough to avoid latency stalls for all but the fastest pixel shaders (which are usually memory bandwidth-bound anyway). There's also the issue of pixel shaders that output to multiple render targets, but that depends on how exactly that feature is implemented. You could run the shader multiple times (not efficient but easiest if you have fixed-size output buffers), or you could run all the render targets through the same ROP (but up to 8 rendertargets with up to 128 bits/pixels – that's a lot of buffer space we're talking), or you could allocate one ROP per output render target.

An of course, if we have these buffers in the ROPs anyway, we might as well treat them as a small cache (i.e. keep them around for a while). This would help if you're drawing lots of small triangles – as long as they're spatially localized, anyway. Again, I'm not sure if GPUs actually do this, but it seems like a reasonable thing to do (you'd probably want to flush these buffers something like once per batch or so though, to avoid the synchronization/coherency issues that full write-back caches bring).

Okay, that explains the memory side of things, and the computational part we've already covered. Next up: Compression!

Depth buffer and color buffer compression

I already explained the basic workings of this in [part 7](https://fgiesen.wordpress.com/2011/07/08/a-trip-through-the-graphics-pipeline-2011-part-7/) (<https://fgiesen.wordpress.com/2011/07/08/a-trip-through-the-graphics-pipeline-2011-part-7/>) while talking about Z; in fact, I don't have much to add about depth buffer compression here. But all the bandwidth issues I mentioned there exist for color values too; it's not so bad for regular rendering (unless the Pixel Shaders output pixels fast enough to hit memory bandwidth limits), but it is a serious issue for MSAA, where we suddenly store somewhere between 2 and 8 samples per pixel. Like Z, we want some lossless compression scheme to save bandwidth in common cases. Unlike Z, plane equations per tile are not a good fit to textured pixel data.

However, that's no problem, because actually, MSAA pixel data is even easier to optimize for: Remember that pixel shaders only run once per pixel, not per sample – unless you're using sample-frequency shading anyway, but that's a D3D11 feature and not commonly used (yet?). Hence, for all pixels that are fully covered by a single primitive, the 2-8 samples stored will usually be the same. And that's the idea behind the common color buffer compression schemes: Write a flag bit (either per pixel, or per quad, or on an even larger granularity) that denotes whether for all the pixels in a compression block, all the per-sample colors are in fact the same. And if that's the case, we only need to store the color once per pixel after all. This is fairly simple to detect during write-back, and again (much like depth compression), it requires some tag bits that we can store in a small on-chip SRAM. If there's an edge crossing the pixels, we need the full bandwidth, but if the triangles aren't too small (and they're basically never *all* small), we can save a good deal of bandwidth on at least part of the frame. And again, we can use the same machinery to accelerate clears.

On the subject of clears and compression, there's another thing to mention: Some GPUs have "hierarchical Z"-like mechanisms that store, for a large block of pixels (a rasterizer tile, maybe even larger) that the block was recently cleared. Then you only need to store one color value for the whole tile (or larger block) in memory. This gives you very fast color clears for some buffers (again, you need some tag bits for this!). However, as soon as any pixel with non-clear color is written to the tile (or larger block), the "this was just cleared" flag needs to be... well, cleared. But we do save a lot of memory bandwidth on the clear itself and the first time a tile is read from memory.

And that's it for our first rendering data path: just Vertex and Pixel Shaders (the most common path). In the next part, I'll talk about Geometry Shaders and how that pipeline looks. But before I conclude this post, I have a small bonus topic that fits into this section.

Aside: Why no fully programmable blend?

Everyone who writes rendering code wonders about this at some point – the regular blend pipeline a serious pain to work with sometimes. So why can't we get fully programmable blend? We have fully programmable shading, after all! Well, we now have the necessary framework to look into this properly. There's two main proposals for this that I've seen – let's look at the both in turn:

1. Blend in Pixel Shader – i.e. Pixel Shader reads framebuffer, computes blend equation, writes new output value.
2. Programmable Blend Unit – "Blend Shaders", with subset of full shader instruction set if necessary. Happen in separate stage after PS.

1. Blend in Pixel Shader

This seems like a no-brainer: after all, we have loads and texture samples in shaders already, right? So why not just allow a read to the current render target? Turns out that unconstrained reads are a *really* bad idea, because it means that every pixel being shaded could (potentially) influence every other pixel being shaded. So what if I reference a pixel in the quad over to the left? Well, a shader for that quad could be running this instant. Or I could be sampling half of my current quad and half of another quads that's currently active – what do I do now? What exactly would be the correct results in that regard, never mind that we'd probably have to shade all quads sequentially to reliably get them? No, that's a can of worms. Unconstrained reads from the frame buffer in Pixel Shaders are out. But what if we get a special render target read instruction that samples one of the active render targets at the current location? Now, that's a lot better – now we only need to worry about writes to the location of the current quad, which is a way more tractable problem.

However, it still introduces ordering constraints; we have to check all quads generated by the rasterizer vs. the quads currently being pixel-shaded. If a quad just generated by the rasterizer wants to write to a sample that'll be written by one of the Pixel Shaders that are currently in flight, we need to wait until that PS is completed before we can dispatch the new quad. This doesn't sound too bad, but how do we track this? We could just have a "this sample is currently being shaded" bit flag... so how many of these bits do we need? At 1920×1080 with 8x MSAA, about 2MB worth of them (that's bytes not bits) – and that memory is global, shared and determines the rate at which we can issue new quads (since we need to mark a quad as busy before we can issue it). Worse, with the hierarchical Z etc. tag bits, they were just a hint; if we ran out of them, we could still render, albeit more slowly. But this memory is *not* optional. We can't guarantee correctness unless we're really tracking every sample! What if we just tracked the "busy" state per pixel (or even quad), and any write to a pixel would block all other such writes? That would work, but it would massively harm our MSAA performance: If we track per sample, we can shade adjacent, non-overlapping triangles in parallel, no problem. But if we track per pixel (or at lower granularity), we effectively serialize all the edge quads. And what happens to our fill rate for e.g. particle systems with lots of overdraw? With the pipeline I described, these render (more or less) as fast as the ROPs can merge the incoming pixels into the store buffers. But if we need to avoid conflicts, we really end up shading the individual overlapping particles in order. This isn't good news for our shader units that are designed to trade latency for throughput, not at all.

Okay, so this whole tracking thing is a problem. What if we just force shading to execute in order? That is, keep the whole thing pipelined and all shaders running in lockstep; now we don't need tracking because pixels will finish in the same order we put them into the pipeline! But the problem here is that we need to make sure the shaders in a batch actually always take the exact same time, which has unfortunate consequences: You always have to wait the worst-case delay time for every texture sample, need to always execute both sides of every branch (someone might at some point need the then/else branches, and we need everything to take the same time!), always runs all loops through for the same number of iterations, can't stop shading on discard... no, that doesn't sound like a winner either.

Okay, time to face the music: Pixel Shader blend in the architecture I've described comes with a bunch of seriously tricky problems. So what about the second approach?

2. "Blend Shaders"

I'll say it right now: This can be made to work, *but...*

Let's just say it has its own problems. For once, we now need another full ALU + instruction decoder/sequencer etc. in the ROPs. This is not a small change – not in design effort, nor in area, nor in power. Second, as I mentioned near the start of this post, our regular "just go wide" tactics don't work so well for blend, because this is a place where we might well get a bunch of quads hitting the same pixels in a row and need to process them in order, so we want low latency. That's a very different design point than our regular unified shader units – so we can't use them for this (it also means texture sampling/memory access in Blend Shaders is a big no, but I doubt that shocks anyone at this point). Third, pure serial execution is out at this point – too low throughput. So we need to pipeline it. But to pipeline it, we need to know how long the pipeline is! For a regular blend unit, it's a fixed length, so it's easy. A blend shader would probably be the same. In fact, due to the design constraints, you're unlikely to get a blend shader – more like a blend register combiner, really, completely with a (presumably relatively low) upper limit on the number of instructions, as determined by the length of the pipeline.

Point being, the serial execution here really constrains us to designs that are still relatively low-level; nowhere near the fully programmable shader units we've come to love. A nicer blend unit with some extra blend modes, you can definitely get; a more open register combiner-style design, possibly, though neither the API guys nor the hardware guys will like it much (the API because it's a fixed function block, the hardware guys because it's big and needs a big ALU+control logic where they'd rather not have it). Fully programmable, with branches, loops, etc. – not going to happen. At that point you might as well bite the bullet and do what it takes to get the "Blend in Pixel Shader" scenario to work properly.

...and that's it for this post! See you next time.

From → Coding, Graphics Pipeline

7 Comments

1. Aras Pranckevičius permalink

First things first: awesome post series!

Now, onto the programmable blending ;) Do you have any information/intuition how some mobile GPUs (PowerVR SGX, NVIDIA Tegra 2, ARM Mali) do programmable blending? They do have it, and not all of them are tile based rasterizers.

Reply

- [fgiesen permalink](#)

Sorry for taking so long to reply, your comment ended up in the spam folder for some reason. Anyway, no, I don't have detailed information, but here's some thoughts:

On all mobile chips, clocks are much lower; memory latency is a bit higher too, but overall, memory wait times are still lower in terms of cycles. That helps. They also have a far lower count of shader units, which means "issue stalls" due to queued quads depending on quads that are still being shaded are less expensive overall (the throughput cost due to a stall is proportional to the number of units that are left idle). Finally, because they have both lower clocks and a less extreme memory:arithmetic latency ratio, they need less quads in-flight per shader to sustain good utilization.

To explain a bit more: The scoreboarding-based scheme I described is one way to avoid "shading races" (and the "most natural" if you're thinking in software terms). Another is to keep track of all in-flight quads. Then every new quad is tested against all in-flight quads for coverage mask collisions (this can be implemented by keeping track of the in-flight quads in a small content-addressable memory, aka CAM). If you have a reasonably small number (say between 16 and 64) quads in-flight at any given time, this works fine. But with 4-16 quads per batch, 10+ batches/Warps/Wavefronts running on a shader unit at any given time, and dozens of shader units, you need an impractically large CAM (and they're power-hungry as hell!). And of course for tile-based renderers the whole scoreboard size issue in the direct scoreboarding algorithm disappears too.

You can make programmable blending work without a huge extra cost if you have fixed latencies (register combiners or shaders without dynamic branching), or if you keep the contested resource has a reasonably small bounded size (e.g. tile-based), or if the number of agents that can conflict is kept small (low number of in-flight quads). But when you have variable latency and impractically large bounds on render target size and number of in-flight quads, you're in trouble. :)

Reply

2. Kevin Rogovin permalink

Wonderful series, I wanted to comment on the blend shader thing in context of mobile.

1) For tiled based renderers, i.e. SGX, Mali, Adreno since the rasterization takes place on (tiny) tile at a time on SRAM, the entire memory pain of a blend shader does not exist. Indeed, there is an ES extension (that Apple iOS now supports) that allows one to read the value of the "framebuffer".

2) Also on mobile, but NOT a tiled based renderer: NVIDIA Tegra (2,3 and 4) also allow one to read the framebuffer value from the fragment shader. One of the icky things is that using the NVIDIA offline compiler, one needs to pragma the blend-state so that it will append those instructions to the fragment shader. I know Tegra is not tiled based, so this design decision I think is odd.

Reply

- o [fgiesen](#) permalink

Yeah, my general discussion applies to both non-tiled and tiled renderers, and I was actually working on a shader compiler for a tiled renderer (with programmable blend inside shaders) at the time I wrote this. :)

The general issue stays the same, though: blending/late Z write is a synchronization point; for any given pixel in the render target, the blend/Z etc. operations have to happen in the right order. There's numerous ways to solve this and all have different trade-offs.

At the very least, the "blend" stages of pixel shading for quads hitting the same location in the render target need to run in the right order, and one way to provide this guarantee (and nothing stronger) is the scoreboarding-like scheme I describe. The problems I discuss wrt render target sizes don't exist in a tiled renderer; all this checking happens tile by tile, so the size is fixed and everything just works. You can also be even more strict and just require that *all* quads blend in the order they were rasterized; this avoids the need for bookkeeping but means that shaders now have a "blend barrier" right before blending starts. This is less bookkeeping but means pixel shader warps/wavefronts are "live" (and potentially stalled) for longer, reducing the available resources for other warps/wavefronts that could actually do useful work instead of just sitting around waiting for their turn at blending. How expensive this is depends on how much it reduces your utilization, which depends on lots of other things including the expected complexity of the shaders you're running.

For a non-tiled renderer like the Tegra, anything that requires explicit per-quad bookkeeping is icky (for the reasons described in the article), which means they're probably using a less precise scheme like the "in-order blend" stuff I described above. Now this doesn't mean that *all* blend operations inside the pipeline have to be synchronized against each other; for example, if you have 4 "shader cores" (the high-level ones, I'm not talking about "CUDA cores" or whatever NVidia's current nomenclature of the day is here), you can just assign a quarter of the render target to each core, usually in some kind of checkerboard pattern. Each shader core "owns" that part of the screen. With that kind of scheme, each shader core only needs to synchronize blending operations against other blending operations it's done by itself, there's no global "locks".

That's one way to make this kind of approach scale; it still suffers once you have complex shaders with several branches and very variable run times though, because one quad that takes long to shade can hold up blending for everything that happens after it.

The ROP design you see in high-end GPUs essentially does the same thing; each ROP owns part of the render target, so they don't conflict and don't need to talk to each other to do the right thing. Blending is still relatively serial work within a ROP. But you can have lots of them, and more importantly, a ROP stalling because some quads aren't done shading yet won't necessarily block the shader units, not until the ROP's input queue fills up anyway. So instead of stalling the shader units (which could be doing all kinds of other work in the mean time), you stall a dedicated unit whose only task is to blend, which has very little state per pixel (much less than the original shader would), and which is designed to be fast enough to "catch up" after most stalls without causing any hitches upstream (in the shader cores, which we want to keep busy).

Reply

Trackbacks & Pingbacks

1. [A trip through the Graphics Pipeline 2011: Index](#) « The ryg blog
2. [A trip through the Graphics Pipeline 2011, part 13](#) « The ryg blog
3. [Photoshop Blend Modes in Unity – The Code Corsair](#)

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 10

July 20, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Welcome back. Last time, we dove into bottom end of the pixel pipeline. This time, we’ll switch back to the middle of the pipeline to look at what is probably the most visible addition that came with D3D10: Geometry Shaders. But first, some more words on how I decompose the graphics pipeline in this series, and how that’s different from the view the APIs will present to you.

There’s multiple pipelines / anatomy of a pipeline stage

This goes back to **part 3** (<https://fgiesen.wordpress.com/2011/07/03/a-trip-through-the-graphics-pipeline-2011-part-3/>), but it’s important enough to repeat it: if you look in, for example, the D3D10 documentation, you’ll find a diagram of the “D3D10 pipeline” that includes all stages that might be active. The “D3D10 pipeline” includes Geometry Shading, even if you don’t have a Geometry shader set, and the same for Stream-Out. In the purely functional model of D3D10, the Geometry Shading stage is always there; if you don’t set a Geometry Shader, it just happens to be very simple (and boring): data is just passed through unmodified to the next pipeline stage(s) (Rasterization/Stream-Out).

That’s the right way to specify the API, but it’s the wrong way to think about it in this series, where we’re concerned with how that functional model is actually implemented in hardware. So how do the two shader stages we’ve seen so far look? For VS, we went through the Input Assembler, which prepared a block of vertices for shading, then dispatched that batch to a shader unit (which chews on it for a while), and then some time later we get the results back, write them into a buffer (for Primitive Assembly), make sure they’re in the right order, then send them down to the next pipeline stage (Culling/Clipping etc.). For PS, we receive to-be-shaded quads from the rasterizer, batch them up, buffer them for a while until a shader unit is free to accept a new batch, dispatch a batch to a shader unit (which chews on it for a while), and then some time later we get the results back, write them into a buffer (for the ROPs), make sure they’re in the right order, then do blend/late Z and send the results on to memory. Sounds kind of familiar, doesn’t it?

In fact, this is how it *always* looks when we want to get something done by the shader units: we need a buffer in the front, then some dispatching logic (which is in fact pretty universal for all shader types and can be shared), then we go wide and run a bunch of shaders in parallel, and finally we need another buffer and a unit that sorts the results (which we received potentially out-of-order from the shader units) back into API order.

We’ve seen shader units (and shader execution) and we’ve seen dispatch; and in fact, now that we’ve seen Pixel Shaders (which have some peculiarities like derivative computation, helper pixels, discard and attribute interpolation), we’re not gonna see any big additions to shader unit functionality until we get to Compute Shaders, with their specialized buffer types and atomics. So for the next few parts, I won’t be talking about the shader units; what’s really different about the various shader types is the shape and interpretation of data that goes in and comes out. The shader parts that don’t deal with IO (arithmetic, texture sampling) stay the same, so I won’t be talking about them.

The Shape of Tris to Shade

So let’s have a look at how our IO buffers for Geometry Shaders look. Let’s start with input. Well, that’s reasonably easy – it’s just what we wrote from the Vertex Shader! Or well, not quite; the Geometry Shader looks at primitives, not individual vertices, so what we really need is the output from Primitive Assembly (PA). Note that there’s multiple ways to deal with this; PA could expand primitives out (duplicating vertices if they’re referenced multiple times), or it could just hand us one block of vertices (I’ll stick with the 32 vertices I used earlier) with an associated small “index buffer” (since we’re indexing into a block of 32 vertices, we just need 5 bits per index). Either way works fine; the former is the natural input format for the clip/cull I discussed after PA, but the latter needs far less buffer space when running GS, so I’ll use that model here.

One reason you need to worry about amount of buffer space with GS is that it can work on some pretty large primitives, because it doesn't just support plain lines or triangles (2 and 3 vertices per primitive respectively), but also lines/triangles with adjacency information (4/6 vertices per primitive). And D3D11 adds input primitives that are much fatter still – a GS can also consume patches with up to 32 control points as input. Duplicating the vertices of e.g. a 16-control point patch, which could each have up to 16 vector attributes (32 with D3D11)? That'd be some serious memory waste. So I'm assuming non-duplicated, indexed vertices for this path. Which makes the input for a batch of primitives: the VS output, plus a (relatively small) index buffer.

Now, the geometry shader runs per primitive. For vertex shaders, we needed to gather a batch of vertices, and we chose our batch size with a simple greedy algorithm that tries to pack as many vertices into a batch as possible without splitting a primitive across multiple batches – fair enough. And for pixel shading, we get plenty of quads from the rasterizer and pack them all into batches. Geometry Shaders are a bit more inconvenient – our input block is guaranteed to contain at least one full primitive, and possibly several – but other than that, the number of primitives in that block completely depends on the vertex cache hit rate. If it's high and we're using triangles, we might get something like 40-43; if we're using triangles with adjacency information we could have as little as 5 if we're unlucky.

Of course, we could try to collect primitives from several input blocks here, but that's kind of awkward too. Now we need to keep multiple input blocks and index buffers around for a single GS batch, and if a single batch can refer to multiple index buffers that means each primitive in that batch now needs to know where to get the indices and vertex data from – more storage requirements, more management, more overhead. Also ugly. And of course even with two input blocks you're still at crappy utilization if you hit two input batches with low vertex cache hit rate. You can support more input blocks, but that eats away at memory – and remember, you need space for the output geometry too (I'll get to that in a bit).

So this is our first snag: with VS, we could basically pick our target batch size, and we chose to not always generate full batches so as to make our lives in PA (and here in the GS, and later in the HS too) a bit easier. With PS, we always shade quads, and even fairly small tris usually hit multiple quads so we get an okay ratio of number of quads to number of tris. But with GS, we don't have full control over either ends of the pipeline (since we're in the middle!), and we need multiple input vertices per primitive (as opposed to multiple quads per one input triangle), so buffering up a lot of input is expensive (both in terms of memory and in the amount of management overhead we get).

At this stage, you can basically pick how many input blocks you're willing to merge to get one block of primitives to geometry shade; that number is going to be low because of the memory requirements (I'd be very surprised to see more than 4), and depending on how important you judge GS to be, you might even pick 1, i.e. don't merge across input blocks at all and live with crappy utilization on GS shading blocks/Warps/Wavefronts! That's not great with triangles and really bad with the primitives that have even more vertices, but not much of an issue when your main use case for GS in practice is expanding points to quads (point sprites) and maybe rendering the occasional cube shadow map (using the Viewport Array Index/Rendertarget Index – I'll get to that in a bit).

GS output: no rose garden over here, either

So how's it looking on the output side? Again, this is more complicated than the plain VS data flow. Much more complicated in fact; while a VS only outputs one thing (shaded vertices) with a 1:1 correspondence between unshaded and shaded vertices, a GS outputs a variable number of vertices (up to a maximum that's specified at compile time), and as of D3D11 it can also have multiple output streams – however, a maximum of one stream can be sent on down the rest of the pipeline, which is the path I'm talking about now. The other destination for GS data (Stream-Out) will be covered in the next part.

A GS produces variable-sized output, but it needs to run with bounded memory requirements (among other things, the amount of memory available for buffers determines how many primitives can be Geometry Shaded in parallel), which is why the maximum number of output vertices is fixed at compile-time. This (together with the number of written output attributes) determines how much buffer space is allocated, and thus indirectly the maximum number of parallel GS invocations; if that number is too low, latency can't be fully hidden, and the GS will stall for some percentage of the time.

Also note that the GS inputs *primitives* (e.g. points, lines, triangles or patches, optionally with adjacency information), but outputs *vertices* – even though we send primitives down to the rasterizer! If the output primitive type is points, this is trivial. For lines and triangles however, we need to reassemble those vertices back into primitives again. This is handled by making the output vertices form a line or triangle strip, respectively. This handles what are perhaps the 3 most important cases well: single lines, triangles, or quads. It's not so convenient if the GS tries to do some actual extrusion or generate otherwise "complicated" geometry, which often needs several "restart strip" markers (which boils down to a single bit per vertex that denotes whether the current strip is continued or a new strip is started). So why the limitation? At the API level, it seems fairly arbitrary – why can't the GS just output a vertex list together with a small index buffer?

The answer boils down to two words: Primitive Assembly. This is what we're doing here – taking a number of vertices and assembling them into a full primitive to send down the pipeline. But we already use that functional block in this data path, just in front of the GS. So for GS, we need a second primitive assembly stage, which we'd like to keep simple, and assembling triangle strips is very simple indeed: a triangle is always 3 vertices from the output buffer in sequential order, with only a bit of glue logic to keep track of the current winding order. In other words, strips are not significantly more complex to support than what is arguably the simplest primitive of all (non-indexed lines/triangles), but they still save output buffer space (and hence give us more potential for parallelism) for typical primitives like quads.

API order again

There's a few problems here, however: in the regular vertex shading path, we know exactly how many primitives there are in a batch and where they are, even before the shaded vertices arrive at the PA buffer – all this is fixed from the point where we set up the batches to shade. If we, for example, have multiple units for cull/clip/triangle setup, they can all start in parallel; they know where to get their vertex data from, and they can know ahead of time which "sequence number" their triangle will have so it can all be put into order.

For GS, we don't generally know how many primitives we're gonna generate before we get the outputs back – in fact, we might not have produced any! But we still need to respect API order: it's first all primitives generated from GS invocation 0, then all primitives from invocation 1, and so on, through to the end of the batch (and of course the batches need to be processed in order too, same as with VS). So for GS, once we get results back, we first need to scan over the output data to determine the locations where complete primitives start. Only then can we start doing cull, clip and triangle setup (potentially in parallel). More extra work!

VPAI and RTAI

These are two features added with GS that don't actually affect Geometry Shader execution, but do have some effect on the processing further downstream, so I thought I'd mention them here: The Viewport Array Index (here, VPAI for short) and Render-target Array Index (RTAI). RTAI first, since it's a bit easier to explain: as you hopefully know, D3D10 adds support for texture arrays. Well, the RTAI gives you render-to-texture-array support: you set a texture array as render target, and then in the GS you can select per-primitive to which array index the primitive should go. Note that because the GS is writing vertices not primitives, we need to pick a single vertex that selects the RTAI (and also VPAI) per primitive; this is always the "leading vertex", i.e. the first specified vertex that belongs to a primitive. One example use case for RTAI is rendering cubemaps in one pass: the GS decides per primitive to which of the cube faces it should be sent (potentially several of them). VPAI is an orthogonal feature which allows you to set multiple viewports and scissor rects (up to 15), and then decide per primitive which viewport to use. This can be used to render multiple cascades in a Cascaded Shadow Map in a single pass, for example, and it can also be combined with RTAI.

As said, both features don't affect GS processing significantly – they're just extra data that gets tacked onto the primitive and then used later: the VPAI gets consumed during the viewport transform, while the RTAI makes it all the way down to the pixel pipeline.

Summary so far

Okay, so there's some amount of trouble on the input end – we don't fully get to pick our input data format, so we need extra buffering on the input data, and even then we have a variable amount of input primitives which we're not necessarily going to be able to partition into nice big batches. And on the output end, we're again assembling a variable number of primitives, don't necessarily know which GS will produce how many primitives in advance (though for some GSs we'll be able to determine this statically from the compiled code, for example because all vertex emits are outside of flow control or inside loops with a known iteration count and no early-outs), and have to spend some time parsing the output before we can send it on to triangle setup.

If that sounds more involved than what we had in the VS-only case, that's because it is. This is why I mentioned above that it's a mistake to think of the GS as something that always runs – even a very simple GS that does nothing except pass the current triangle through goes through two more buffering stages, an extra round of primitive assembly, and might execute on the shader units with poor utilization. All of this has a cost, and it tends to add up: I checked it when D3D10 hardware was fairly new, and on both AMD and NVidia hardware, even a pure pass-through GS was between 3x and 7x slower than no GS at all (in a geometry-

limited scenario, that is). I haven't re-run this experiment on more recent hardware; I would assume that it's gotten better by now (this was the first generation to implement GS, and features don't usually have good performance in the first GPU generation that implements them), but the point still stands: just sending something through the GS pipe, even if nothing at all happens there, has a very visible cost.

And it doesn't help that GSs produce primitives as strips, sequentially; for a Vertex Shader, we get one invocation per vertex, which reads one vertex and writes one vertex (nice). For a GS, though, we might end up having only a batch of 11 GSs running (because there wasn't enough primitives in the input buffer), with each of them running fairly long and producing something like 8 output vertices. That's a long time to be running at low utilization! (Remember we need somewhere between 16 and 64 independent jobs per batch we dispatch to the shader units). It's even more annoying if the GS mainly consists of a loop – for example, in the "render to cube map" case I mentioned for RTAI, we loop over the 6 faces in a cube, check if a triangle is visible on that face, and output a triangle if that's the case. The computations for the 6 faces are really independent; if possible, we'd like to run them in parallel!

Bonus: GS Instancing

Well, enter GS Instancing, another feature new in D3D11 – poorly documented, sadly (and I'm not sure if there's any good examples for it in the SDK). It's fairly simple to explain, though: for each input primitive, the GS gets run not just once but multiple times (this is a static count selected at compile time). It's basically equivalent to wrapping the whole shader in a

```
for (int i = 0; i < N; i++)
{
    // ...
}
```

block, only the loop is handled outside the shader by actually generating multiple GS invocations per input primitive, which helps us get larger batch sizes and thus better utilization. The `i` is exported to the shader as a system-generated value (in D3D11, with Semantic `SV_GSInstanceID`). So if you have a GS like that, just get rid of the outer loop, add a `[instances(N)]` declaration and declare `i` as input with the right semantic and it'll probably run faster for very little work on your part – the magic of giving more independent jobs to a massively parallel machine!

Anyway, that's it on Geometry Shaders. I've skipped Stream-Out, but this post is already long enough, and besides SO is a big enough topic (and independent enough of GS!) to warrant its own post. Next post, to be more precise. Until then!

From → Coding, Graphics Pipeline

2 Comments

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog
2. Page not found « The ryg blog

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 11

August 14, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Welcome back! This time, the focus is going to be on Stream-Out (SO). This is a facility for storing the Output of the Geometry Shader stage to memory, instead of sending it down the rest of the pipeline. This can be used to e.g. cache skinned vertex data, or as a sort of poor man’s Compute Shader on D3D10-level hardware using the D3D10 API (note that with D3D11, you can just use CS 4.0, even on D3D10 hardware). And just like the GS Instancing I mentioned last time, some of this is very poorly described in the API docs, so I’ll have a few comments about API usage even though it’s technically out of the intended scope of this series.

Vertex Shader Stream-Out (i.e. SO with NULL GS)

This is one of the features that’s not properly explained in the D3D10 (or D3D11, for that matter) docs; in fact, it’s not mentioned there at all except for a small throwaway remark in “Getting Started with the Stream-Output Stage (Direct3D 10)”. You’re supposed to figure it out from the examples – which themselves don’t exactly go out of their way to make it clear what’s going on. That’s a pity – VS Stream-Out is easier than GS SO, and has some pretty useful applications by itself (e.g. caching skinned vertices).

So here’s how it’s done in D3D10 and 11: You simply pass Vertex Shader bytecode (instead of GS bytecode) to `CreateGeometryShaderWithStreamOutput`. Yes, the docs mention something about “Size of the compiled geometry shader” here – ignore it. What you get back is a Geometry Shader object that you can then pass to `GSSetShader`. This is, in effect, a NULL Geometry Shader – it doesn’t actually go through GS processing. It’s just some wrapper (more like duct tape really) to make it fit into the API model, where all rendering passes through the GS stage and SO comes right after GS – though as I’ve explained last time, actual HW tends to skip the GS stage completely when there’s no GS set.

So the shaded vertices get assembled into primitives as before, but instead of getting sent down the rest of the pipeline as already described, they get forwarded to Stream-Out, where they arrive – as always – in a buffer. What exactly happens with them then depends on the Stream-Out declaration (which is passed at creation time). In the Stream-Out declaration, the app gets to specify where it wants each output vector to end up in the Stream-Out targets (or SO targets for short). If the SO declaration “matches” the Vertex Shader Output Declaration (i.e. the same attributes in the same order), data from the input buffers can be streamed more or less unprocessed into memory. If it doesn’t match the declaration exactly – it might skip some attributes written by the shader, or write them in a different order – either way, there’s some extra reordering involved. This might involve a dedicated reordering unit (which basically implements a gather-type operation from the SO input buffers), or it might involve generating lots of small memory writes instead of large burst writes, or something similar. Either way, it’s extra effort and generally slower; the details of what exactly triggers a slow path depend on the hardware specifics, but really, it doesn’t matter that much. If you want optimal SO performance, just make sure the SO declaration and Output declarations agree.

Another point is that SO usually doesn’t have access to a very high-performance path to the memory subsystem. Unlike e.g. the ROPs, SO isn’t really (yet?) a full citizen in current GPU designs, so it often only has access to one memory channel or something of the sort. That’s something to keep in mind if you’re producing a lot of data via SO. This is compounded by SO outputs always being full floats, so there’s no way to conserve bandwidth by using one of the packed vertex data types.

Final remark on VS SO: As I mentioned earlier, SO operates on assembled *primitives*, not individual vertices. Note that Primitive Assembly discards adjacency information if it makes it that far down the pipeline, and since this happens before SO, vertices corresponding to adjacency info won’t make it into SO buffers either. SO working on primitives not individual vertices is relevant for use cases like instancing a single skinned mesh (in a single pose) several times. If you were to draw your triangle mesh as you usually would and then use SO on that, this results in a data explosion – you get 3 unpacked, unshared vertices per input primitive. This works, but isn’t exactly an efficient use of bandwidth, both on the SO and the later vertex input side. Instead, you should draw your triangle mesh as a (non-indexed) point list in the first pass, thereby shading each vertex exactly once. The SO buffer then ends up in 1:1 correspondence to your original vertex buffer, only with skinned instead of non-skinned vertices. You can then use that vertex buffer with your original primitive topology and index buffer.

Geometry Shader SO: Multiple streams

This basically works like SO with a NULL GS, except there's a Geometry Shader involved, which adds some new capabilities (and complications). In the VS case, we just had one output stream (note that streams are a D3D11+ feature – they don't exist on D3D10-level HW). That stream could be sent to SO or not, and it could also be sent to down the pipeline to viewport/clip/cull or not, but that's it. But Geometry Shaders allow multiple streams, which makes output routing a bit more difficult.

Basically, every GS can write to (as of D3D11) up to 4 streams. Each stream may be sent on to SO targets – yes, plural: a single stream can write to multiple SO targets, but a single SO target can receive values from only one stream, i.e. this is a one-to-many relationship, not a fully general many-to-many one. The presence of streams has some implications for SO buffering – instead of a single input buffer like I described in the NULL GS case, we now may have multiple input buffers, one per stream. In addition to SO targets, up to one stream may be sent down the pipe – i.e. the regular rendering pipeline and SO may be used simultaneously.

As in the NULL GS case, SO works on primitives, not individual vertices – that is, the strips you output in the GS get expanded out to full lines or triangles before they get into SO.

Tracking output size

There's another issue here: we don't necessarily know how much output data is going to be produced from SO. For GS, this comes about because each GS invocation may produce a variable number of output primitives; but even in the simpler VS case, as soon as indexed primitives are involved, the app might slip some "primitive cut" indices in there that influence how many primitives actually get written. This is a problem if we then want to draw from that SO buffer later, because we don't know how many vertices are actually in there! We do have an upper bound – the maximum capacity of the buffer as created – but that's it. Now, this could be resolved using some kind of query mechanism, but once you think it through, that seems fairly backwards: at the point we're using the SO buffer for drawing, we obviously do know how many primitives we actually wrote – the SO unit needs to keep track of its current output position, after all! If we employed some query mechanism, we would end up transporting that single 32-bit value back over the bus to the driver, which passes it on to the API, which passes it on to the app – which then immediately dispatches another draw, going through all the layers again in the opposite direction.

So that's not how it's solved. Instead, there's `DrawAuto`. The idea is very simple – the GPU already knows how many valid vertices it actually wrote to the output buffer; the SO unit keeps track of that while it's writing, and the final counter is also kept in memory (along with the buffer) since the app may render to a SO buffer in multiple passes. This counter is then used for `DrawAuto`, instead of having the app submit an explicit count itself – simplifying things considerably and avoiding the costly round-trip completely. Note that this query mechanism does exist – both for checking the number of vertices written and to determine whether an overflow occurred. But it's not on the critical path for rendering from SO buffers, which makes things a lot simpler for driver developers.

And that's it for SO, really. Not really a lot of HW info in this one, and not really a super-interesting topic from a pipeline perspective, which is why it took me so long to finish; sorry about that. Next up is Tessellation – this should be a lot quicker, since it's a fun topic :)

From → Coding, Graphics Pipeline

One Comment

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 12

September 6, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Welcome back! This time, we’ll look into what is perhaps the “poster boy” feature introduced with the D3D11 / Shader 5.x hardware generation: Tessellation. This one is interesting both because it’s a fun topic, and because it marks the first time in a long while that a significant user-visible component has been added to the graphics pipeline that’s not programmable.

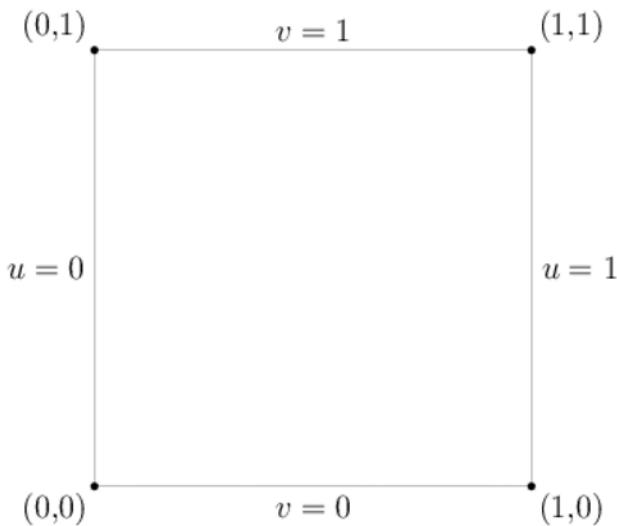
Unlike Geometry Shaders, which are conceptually quite easy (it’s just a shader that sees whole primitives as opposed to individual vertices), the topic of “Tessellation” requires some more explanation. There’s tons of ways to tessellate geometry – to name just the most popular ones, there’s Spline Patches in dozens of flavors, various types of Subdivision Surfaces, and Displacement Mapping – so from the bullet point “Tessellation” alone it’s not at all obvious what services the GPU provides us with, and how they are implemented.

To describe how hardware tessellation works, it’s probably easiest to start in the middle – with the actual primitive tessellation step, and the various requirements that apply to it. I’ll get to the new shader types (Hull Shaders and Domain Shaders in D3D11 parlance, Tessellation Control Shader and Tessellation Evaluation Shader in OpenGL 4.0 lingo) later.

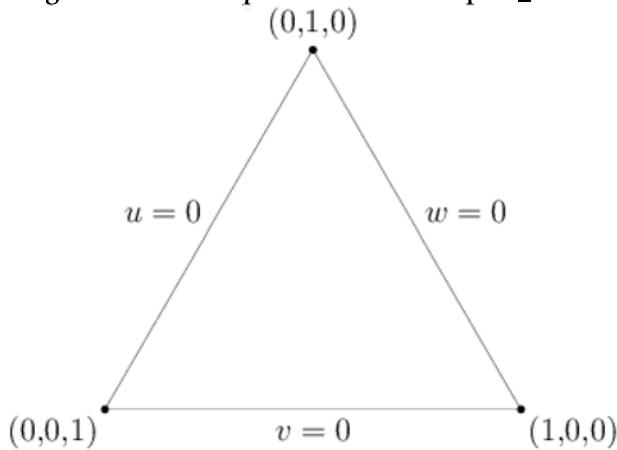
Tessellation – not quite like you’d expect

Tessellation as implemented by Shader 5.x class HW is of the “patch-based” variety. Patch types in the CG literature are mostly named by what kind of function is used to construct the tessellated points from the control points (B-spline patches, Bézier triangles, etc.). But we’ll ignore that part for now, since it’s handled in the new shader types. The actual fixed-function tessellation unit deals only with the *topology* of the output mesh (i.e. how many vertices there are and how they’re connected to each other); and it turns out that from this perspective, there’s basically only two different types of patches: quad-based patches, which are defined on a parameter domain with two orthogonal coordinate axes (which I’ll call u and v here, both are in $[0,1]$) and usually constructed as a tensor product of two one-parameter basis functions, and triangle-based patches, which use a redundant representation with three coordinates (u, v, w) based on barycentric coordinates (i.e. $u, v, w \geq 0, u + v + w = 1$). In D3D11 parlance, these are the “quad” and “tri” domains, respectively. There’s also an “isoline” domain which instead of a 2D surface produces one or multiple 1D curves; I’ll treat it the same way as I did lines and point primitives throughout this series: I acknowledge its existence but won’t go into further detail.

Tessellated primitives can be drawn naturally in their respective domain coordinate systems. For quads, the obvious choice of drawing the domain is as a unit square, so that’s what I’ll use; for triangles, I’ll use an equilateral triangle to visualize things. Here’s the coordinate systems I’ll be using in this post with both the vertices and edges labeled:

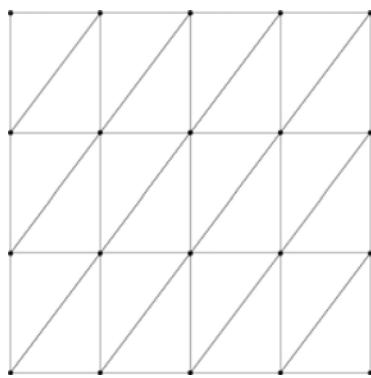


(https://fgiesen.files.wordpress.com/2011/09/quad_coords2.png)

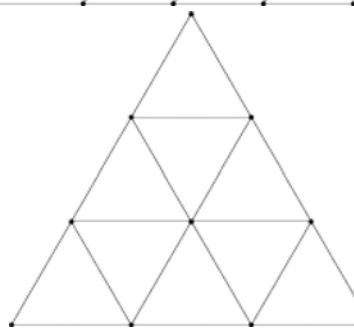


(https://fgiesen.files.wordpress.com/2011/09/tri_coords.png)

Anyway, both triangles and quads have what I would consider a “natural” way to tessellate them, depicted below. But it turns out that’s not actually the mesh topology you get.

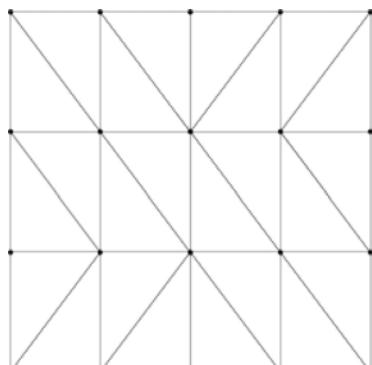


(https://fgiesen.files.wordpress.com/2011/09/quad_tess_simple1.png)

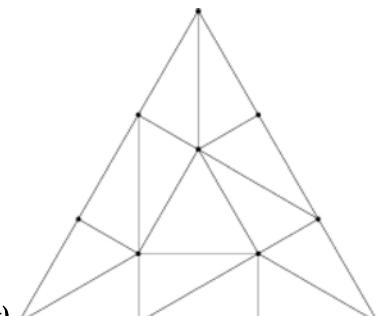


(https://fgiesen.files.wordpress.com/2011/09/tri_tess_simple1.png)

Here’s the *actual* meshes that the tessellator will produce for the given input parameters:



(https://fgiesen.files.wordpress.com/2011/09/quad_tess4x3.png)
https://fgiesen.files.wordpress.com/2011/09/tri_tess3.png)



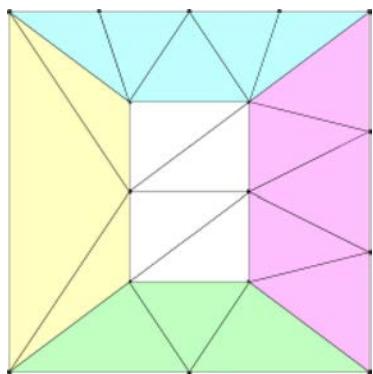
For quads, this is (roughly) what we're expecting – except for some flipped diagonals, which I'll get to in a minute. But the triangle is a completely different beast. It's got a very different topology from the "natural" tessellation I showed above, including a different number of vertices (12 instead of 10). Clearly, there's something funny going on here – and that something happens to be related to the way transitions between different tessellation levels are handled.

Making ends meet

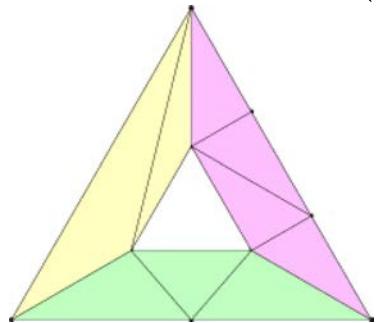
The elephant in the room is handling transitions between patches. Tessellating a single triangle (or quad) is easy, but we want to be able to determine tessellation factors per-patch, because we only want to spend triangles where we need them – and not waste tons of triangles on some distant (and possibly backface-culled) parts of the mesh. Additionally, we want to be able to do this quickly and ideally without extra memory usage; that means a global fix-up post-pass or something of that caliber is out of the question.

The solution – which you've already encountered if you've written a Hull or Domain shader – is to make all of the actual tessellation work purely local and push the burden of ensuring watertightness for the resulting mesh down to the shaders. This is a topic all by itself and requires, among other things, **great care in the Domain Shader code** (<http://www.ludicon.com/castano/blog/2010/09/precise/>); I'll skip all the details about expression evaluation in shaders and stick with the basics. The basic mechanism is that each patch has multiple tessellation factors (TFs), which are computed in the Hull Shader: one or two for the actual inside of the patch, plus one for each edge. The TFs for the inside of the patch can be chosen freely; but if two patches share an edge, they'd better compute the exact same TFs along that edge, or there will be cracks. The hardware doesn't care – it will process each patch by itself. If you do everything correctly, you'll get a nice watertight mesh, otherwise – well, that's your problem. All the HW needs to make sure is that it's *possible* to get watertight meshes, preferably with reasonable efficiency. That by itself turns out to be tricky in some places; I'll get to that later.

So, here are some new reference patches – this time with different TFs along each edge so we can see how that works:



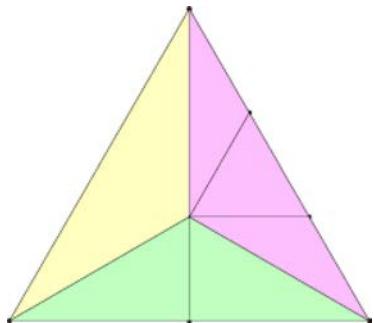
(https://fgiesen.files.wordpress.com/2011/09/quad_tess_asym.png)



(https://fgiesen.files.wordpress.com/2011/09/tri_tess_asym.png)

I've colored the areas influenced by the different edge tessellation factors; the uncolored center part in the middle only depends on the inside TFs. In these images, the $u=0$ (yellow) edge has a TF of 2, the $v=0$ (green) edge has a TF of 3, the $u=1 / w=0$ (pink) edge has a TF of 4, and the $v=1$ (quad only, cyan) edge has a TF of 5 – exactly the number of vertices along the corresponding outer edge. As should be obvious from these two images, the basic building block for topology is just a nice way to stitch two subdivided edges with different number of vertices to each other. The details of this are somewhat tricky, but not particularly interesting, so I won't go into it.

As for the inside TFs, quads are fairly easy: The quad above has an inside TF of 3 along u and 4 along v . The geometry is basically that of a regular grid of that size, except with the first and last rows/columns replaced by the respective stitching triangles (if any edge has a TF of 1, the resulting mesh will have the same structure as if the inside TFs for u/v were both 2, even if they're smaller than that). Triangles are a bit more complicated. Odd TFs we've already seen – for a TF of $\frac{N+1}{2}$, they produce a mesh consisting of $\frac{N+1}{2}$ concentric rings, the innermost of which is a single triangle. For even TFs, we get $\frac{N}{2}$ concentric rings with a center vertex instead of a center triangle. Below is an image of the simplest even case, $N=2$, which consists just of edge stitches plus the center vertex.



(https://fgiesen.files.wordpress.com/2011/09/tri_tess_asym_even.png)

Finally, when triangulating quads, the diagonal is generally chosen to point away from the center of the patch (in the domain coordinate space), with a consistent tie-breaking rule. This is simply to ensure maximum rotational symmetry of the resulting meshes – if there's extra degrees of freedom, might as well use them!

Fractional tessellation factors and overall pipeline flow

So far, I've only talked about integer TFs. In two of the so-called “partitioning types”, namely “Integer” and “Pow2”, that's all the Tessellator sees. If the shader generates a non-integer (or, respectively, non-power-of-2) TF, it will simply get rounded up to the next acceptable value. More interesting are the remaining two partitioning types: Fractional-odd and Fractional-even tessellation.

Instead of jumping from tessellation factor to tessellation factor (which would cause visible pops), new vertices start out at the same position as an existing vertex in the mesh and then gradually move to their new positions as the TF increases.

For example, with fractional-odd tessellation, if you were to use an inner TF of 3.001 for the above triangle, the resulting mesh would look very much like the mesh for a TF of 3 – but topologically, it'd be the same as if the TF was 5, i.e. it's a patch with 3 concentric rings, even though the middle ring is very narrow. Then as the TF gets closer to 5, the middle ring expands until it is eventually at its final position for TF 5. Once you raise the TF past 5, the mesh will be topologically the same as is the TF was 7, but again with a number of almost-degenerate triangles in the middle, and so forth. Fractional-even tessellation uses the same principle, just with even TFs.

The output of the tessellator then consists of two things: First, the positions of the tessellated vertices in domain coordinates, and second, the corresponding connectivity information – basically an index buffer.

Now, with the basic function of the fixed-function tessellator unit explained, let's step back and see what we need to do to actually churn out primitives: First, we need to input a bunch of input control points comprising a patch into the Hull Shader. The HS then computes output control points and "patch constants" (both of which get passed down to the Domain Shader), plus all Tessellation Factors (which are essentially just more patch constants). Then we run the fixed-function tessellator, which gives us a bunch of Domain Positions to run the Domain Shader at, plus the associated indices. After we've run the DS, we then do another round of primitive assembly, and then send the primitives either down to the GS pipeline (if it's active) or Viewport transform, Clip and Cull (if not).

So let's look a bit into the HS stage.

Hull Shader execution

Like **Geometry Shaders** (<https://fgiesen.wordpress.com/2011/07/20/a-trip-through-the-graphics-pipeline-2011-part-10/>), Hull Shaders work on full (patch) primitives as input – with all the input buffering headaches that causes. How much of a headache entirely depends on the type of input patch. If the patch type is something like a cubic Bézier patch, we need $4 \times 4 = 16$ input points *per patch* and might just produce a single quad of output (or even none at all, if the patch is culled); clearly, that's a somewhat awkward amount of data to work with, and doesn't lend itself to very efficient shading. On the other hand, if tessellation takes plain triangles as input (which a lot of people do), input buffering is pretty tame and not likely to be a source of problems or bottlenecks.

More importantly, unlike Geometry Shaders (which run for every primitive), Hull Shaders don't run all that often – they run once *per patch*, and as long as there's any actual tessellation going on (even at modest TFs), we have way less patches than we have output triangles. In other words, even when HS input is somewhat inefficient, it's less of an issue than in the GS case simply because we don't hit it that often.

The other nice attribute of Hull Shaders is that, unlike Geometry Shaders, they don't have a variable amount of output data; they produce a fixed amount of control points, each with a fixed amount of associated attributes, plus a fixed amount of patch constants. All of this is statically known at compile time; no dynamic run-time buffer management necessary. If we Hull Shade 16 hulls at a time, we know exactly where the data for each hull will end up before we even start executing the shader. That's definitely an advantage over Geometry Shaders; for lots of Geometry Shaders, it's possible to know statically how many output vertices will be generated (for example because all the control flow leading to `emit` / `cut` instructions can be statically evaluated at compile time), and for all of them, there's a guaranteed maximum number of output vertices, but for HS, we have a guaranteed fixed amount of output data, no additional analysis required. In short, there's no problems with output buffer management, other than the fact that, again depending on the primitive type, we might need lots of output buffer space which limits the amount of parallelism we can achieve (due to memory/register constraints).

Finally, Hull Shaders are somewhat special in the way they are compiled in D3D11; all other shader types basically consist of one block of code (with some subroutines maybe), but Hull Shaders are generated factored into multiple phases, each of which can consist of multiple (independent) threads of execution. The details are mainly of interest to driver and shader compiler programmers, but suffice it to say that your average HS comes packaged in a form that exposes lots of latent parallelism, if there is any. It certainly seems like Microsoft was really keen to avoid the bottlenecks that plague Geometry Shaders this time around.

Anyway, Hull Shaders produce a bunch of output per patch; most of it is just kept around until the corresponding Domain Shaders run, except for the TFs, which get sent to the tessellator unit. If any of the TFs are less than or equal to zero (or NaN), the patch is culled, and the corresponding control points and patch constants silently get thrown away. Otherwise, the Tessellator (which implements the functionality described above) kicks in, reads the just-shaded patches, and starts churning out domain point positions and triangle indices, and we need to get ready for DS execution.

Domain Shaders

Just like for Vertex Shading (<https://fgiesen.wordpress.com/2011/07/03/a-trip-through-the-graphics-pipeline-2011-part-3/>) way back, we want to gather multiple domain vertices into one batch that we shade together and then pass on the PA. The fixed-function tessellator can take care of this: “just” handle it along with producing vertex positions and indices (I put the “just” in quotes here because this does involve some amount of bookkeeping).

In terms of input and output, Domain Shaders are very simple indeed: the only input they get that actually varies per vertex is the domain point u and v coordinates (w , when used, doesn’t need to be computed or passed in by the tessellator; since $u + v + w = 1$, it can be computed as $w = 1 - u - v$). Everything else is either patch constants, control points (all of which are the same across a patch) or constant buffers. And output is basically the same as for Vertex Shaders.

In short, once we get to the DS, life is good; the data flow is almost as simple as for VS, which is a path we know how to run efficiently. This is perhaps the biggest advantage of the D3D11 tessellation pipeline over Geometry Shaders: the actual triangle amplification doesn’t happen in a shader, where we waste precious ALU cycles and need to keep buffer space for a worst-case estimate of vertices, but in a localized element (the tessellator) that is basically a state machine, gets very little input (a few TFs) and produces very compact output (effectively an index buffer, plus a 2D coordinate per output vertex). Because of this, we need way less memory for buffering, and can keep our Shader Units busy with actual shading work instead of housekeeping.

And that’s it for this post – next up: Compute Shaders, aka the final part in my original outline for this series! Until then.

Final remarks

As usual, I cut a few corners. There’s the “isoline” patch type, which I didn’t go into at all (if there’s any demand for this, I can write it up). The Tessellator has all kinds of symmetry and precision requirements; as far as vertex domain positions are concerned, you can basically expect bit-exact results between the different HW vendors, because the D3D11 spec really nails this bit down. What’s intentionally not nailed down is the order in which vertices or triangles are produced – an implementation can do what it wants there, provided it does so consistently (i.e. the same input has to produce the same output, always). There’s a bunch of subtle constraints that go into this too – for example, all domain positions written by the Tessellator need to have both u and $1-u$ (and also v and $1-v$) exactly representable as float; there’s a bunch of necessary conditions like this so that Domain Shaders can then produce watertight meshes (this rule in particular is important so that a shared edge AB between two patches, which is AB to one patch and BA to the other, can get tessellated the same way for both patches).

Writing Domain Shaders so they actually can’t produce cracks is tricky and requires great care; I intentionally sidestep the topic because it’s outside the scope of this series. Another much more trivial issue that I didn’t mention is the winding order of triangles generated by the Tessellator (answer: it’s up to the App – both clockwise and counterclockwise are supported).

The description of Input/Output buffering for Hull and Domain shaders is somewhat terse, but it’s very similar to stages we’ve already seen, so I’d rather keep it short and avoid extra clutter; re-read the posts on Vertex Shaders and Geometry Shaders if this was too fast.

Finally, because the Tessellation pipeline can feed into the GS, there’s the question of whether it can generate adjacency information. For the “inside” of patches this would be conceivable (just more indices for the Tessellator unit to write), but it gets ugly fast once you reach patch edges, since cross-patch adjacency needs exactly the kind of global “mesh awareness” that the Tessellation pipeline design tries so hard to avoid. So, long story short, no, the tessellator will not produce adjacency information for the GS, just plain triangles.

From → Coding, Graphics Pipeline

12 Comments

1. Naoki permalink

Awesome article, I learnt a lot. I have a question though. Could you elaborate a bit on how you compute inside tessellation factors? For example in the image of the square above, why is the inside TF equal to 3 along u and 4 along v ?

Thanks.

Reply

- [fgiesen permalink](#)

D3D inside TFs for quad domains count all rows and columns of the grid, including the ones covered by edge regions. If the inside TF were less than 2 (but still >0) in either u or v, the white “inside” region would just disappear completely. I mention this in the article – the “(if any edge has a TF of 1, the resulting mesh will have the same structure as if the inside TFs for u/v were both 2, even if they’re smaller than that)” bit. By far the easiest way to understand how the factors interact is to write a small app that renders a single primitive in wireframe mode and then play with the values for a bit. :)

Reply

- [Naoki permalink](#)

Yes, I'll probably do that :) Thanks for the explanation.

2. [Thanh Nguyen permalink](#)

Another awesome article! Your “trip through the graphics pipeline” series is very informative and helpful. I have a question though: your image for quad primitive vertices/edges ordering is in ccw order. D3D11 doc is saying they should be in cw order: [http://msdn.microsoft.com/en-us/library/windows/desktop/ff471574\(v=vs.85\).aspx](http://msdn.microsoft.com/en-us/library/windows/desktop/ff471574(v=vs.85).aspx)

So which one is the correct order?

Reply

- [fgiesen permalink](#)

Sorry for taking a while to reply. I show the coordinate system I use at the top of the article. I use the standard mathematical convention with (0,0) being the origin, the positive x axis pointing right and the positive y axis pointing up. In that coordinate system, the D3D edge ordering (u=0, then v=0, then u=1, then v=1) is a counter-clockwise sweep. The D3D docs put the v=0 axis at the “top of the patch”, which means they refer to a patch coordinate system where (0,0) is in the top-left corner and positive y points downwards. That corresponds to a y-flip of the coordinate system, which turns counter-clockwise sweeps into clockwise sweeps (and vice versa). If you prefer to draw it that way, just mirror all my figures about the x axis :)

It's fairly arbitrary either way, since the patch uv coordinate system is something that's entirely up to the user; you get to pick whatever you like. The only thing you need to be careful about is the winding order of the output triangles (the “output topology”): D3D supports both “triangle_cw” and “triangle_ccw”, but they refer to clockwise or counterclockwise in their UV space with (0,0) being the top left corner.

Reply

- [Thanh Nguyen permalink](#)

Ah, I see. Thanks for the clarification.

3. [Sven Kiesser permalink](#)

Great article and well explained, i have one question though:

Why are there two different types of fractional_spacing (odd and even)? Is one type not enough?

Reply

- [fgiesen permalink](#)

I don't know what the rationale was for including both. For what it's worth, fractional_odd seems to be far more popular in practice.

Reply

4. [Paul Frischknecht permalink](#)

Thanks. The sample “AdaptiveTessellationCS40” that was included with some dx sdks shows the formulas/algorithms for the tessellation strategies for those interested.

Reply

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog
2. » Quad Patch Tessellation in Unity Defective Studios Devblog
3. Good resource on tessellation subdivison pattern | Technology & Programming Answers

Blog at WordPress.com.

The ryg blog

When I grow up I'll be an inventor.

A trip through the Graphics Pipeline 2011, part 13

October 9, 2011

This post is part of the series “*A trip through the Graphics Pipeline 2011*” (<https://fgiesen.wordpress.com/2011/07/09/a-trip-through-the-graphics-pipeline-2011-index/>).

Welcome back to what’s going to be the last “official” part of this series – I’ll do more GPU-related posts in the future, but this series is long enough already. We’ve been touring all the regular parts of the graphics pipeline, down to different levels of detail. Which leaves one major new feature introduced in DX11 out: Compute Shaders. So that’s gonna be my topic this time around.

Execution environment

For this series, the emphasis has been on overall dataflow at the architectural level, not shader execution (which is explained well elsewhere). For the stages so far, that meant focusing on the input piped into and output produced by each stage; the way the internals work was usually dictated by the shape of the data. Compute shaders are different – they’re running by themselves, not as part of the graphics pipeline, so the surface area of their interface is much smaller.

In fact, on the input side, there’s not really any buffers for input data at all. The only input Compute Shaders get, aside from API state such as the bound Constant Buffers and resources, is their thread index. There’s a tremendous potential for confusion here, so here’s the most important thing to keep in mind: a “thread” is the atomic unit of dispatch in the CS environment, and it’s a substantially different beast from the threads provided by the OS that you probably associate with the term. CS threads have their own identity and registers, but they don’t have their own Program Counter (Instruction Pointer) or stack, nor are they scheduled individually.

In fact, “threads” in CS take the place that individual vertices had during **Vertex Shading** (<https://fgiesen.wordpress.com/2011/07/03/a-trip-through-the-graphics-pipeline-2011-part-3/>), or individual pixels during **Pixel Shading** (<https://fgiesen.wordpress.com/2011/07/10/a-trip-through-the-graphics-pipeline-2011-part-8/>). And they get treated the same way: assemble a bunch of them (usually, somewhere between 16 and 64) into a “Warp” or “Wavefront” and let them run the same code in lockstep. CS threads don’t get scheduled – Warps and Wavefronts do (I’ll stick with “Warp” for the rest of this article; mentally substitute “Wavefront” for AMD). To hide latency, we don’t switch to a different “thread” (in CS parlance), but to a different Warp, i.e. a different bundle of threads. Single threads inside a Warp can’t take branches individually; if at least one thread in such a bundle wants to execute a certain piece of code, it gets processed by all the threads in the bundle – even if most threads then end up throwing the results away. In short, CS “threads” are more like SIMD lanes than like the threads you see elsewhere in programming; keep that in mind.

That explains the “thread” and “warp” levels. Above that is the “thread group” level, which deals with – who would’ve thought? – groups of threads. The size of a thread group is specified during shader compilation. In DX11, a thread group can contain anywhere between 1 and 1024 threads, and the thread group size is specified not as a single number but as a 3-tuple giving thread x, y, and z coordinates. This numbering scheme is mostly for the convenience of shader code that addresses 2D or 3D resources, though it also allows for traversal optimizations. At the macro level, CS execution is dispatched in multiples of thread groups; thread group IDs in D3D11 again use 3D group IDs, same as thread IDs, and for pretty much the same reasons.

Thread IDs – which can be passed in in various forms, depending on what the shader prefers – are the only input to Compute Shaders that’s not the same for all threads; quite different from the other shader types we’ve seen before. This is just the tip of the iceberg, though.

Thread Groups

The above description makes it sound like thread groups are a fairly arbitrary middle level in this hierarchy. However, there's one important bit missing that makes thread groups very special indeed: Thread Group Shared Memory (TGSM). On DX11 level hardware, compute shaders have access to 32k of TGSM, which is basically a scratchpad for communication between threads in the same group. This is the primary (and fastest) way by which different CS threads can communicate.

So how is this implemented in hardware? It's quite simple: all threads (well, Warps really) within a thread group get executed by the same shader unit. The shader unit then simply has at least 32k (usually a bit more) of local memory. And because all grouped threads share the same shader unit (and hence the same set of ALUs etc.), there's no need to include complicated arbitration or synchronization mechanisms for shared memory access: only one Warp can access memory in any given cycle, because only one Warp gets to issue instructions in any cycle! Now, of course this process will usually be pipelined, but that doesn't change the basic invariant: per shader unit, we have exactly one piece of TGSM; accessing TGSM might require multiple pipeline stages, but actual reads from (or writes to) TGSM will only happen inside one pipeline stage, and the memory accesses during that cycle all come from within the same Warp.

However, this is not yet enough for actual shared-memory communication. The problem is simple: The above invariant guarantees that there's only one set of accesses to TGSM per cycle even when we don't add any interlocks to prevent concurrent access. This is nice since it makes the hardware simpler and faster. It does not guarantee that memory accesses happen in any particular order from the perspective of the shader program, however, since Warps can be scheduled more or less randomly; it all depends on who is runnable (not waiting for memory access / texture read completion) at certain points in time. Somewhat more subtle, precisely because the whole process is pipelined, it might take some cycles for writes to TGSM to become "visible" to reads; this happens when the actual read and write operations to TGSM occur in different pipeline stages (or different phases of the same stage). So we still need some kind of synchronization mechanism. Enter barriers. There's different types of barriers, but they're composed of just three fundamental components:

1. *Group Synchronization*. A Group Synchronization Barrier forces all threads inside the current group to reach the barrier before any of them may consume past it. Once a Warp reaches such a barrier, it will be flagged as non-runnable, same as if it was waiting for a memory or texture access to complete. Once the last Warp reaches the barrier, the remaining Warps will be reactivated. This all happens at the Warp scheduling level; it adds additional scheduling constraints, which may cause stalls, but there's no need for atomic memory transactions or anything like that; other than lost utilization at the micro level, this is a reasonably cheap operation.
2. *Group Memory Barriers*. Since all threads within a group run on the same shader unit, this basically amounts to a pipeline flush, to ensure that all pending shared memory operations are completed. There's no need to synchronize with resources external to the current shader unit, which means it's again reasonably cheap.
3. *Device Memory Barriers*. This blocks all threads within a group until all memory accesses have completed – either direct or indirect (e.g. via texture samples). As explained earlier in this series, memory accesses and texture samples on GPUs have long latencies – think more than 600, and often above 1000 cycles – so this kind of barrier will really hurt.

DX11 offers different types of barriers that combine several of the above components into one atomic unit; the semantics should be obvious.

Unordered Access Views

We've now dealt with CS input and learned a bit about CS execution. But where do we put our output data? The answer has the unwieldy name "unordered access views", or UAVs for short. An UAV seems somewhat similar to render targets in Pixel Shaders (and UAVs can in fact be used in addition to render targets in Pixel Shaders), but there's some very important semantic differences:

Most importantly, as the name suggests, access to UAVs is "unordered", in the sense that the API does not guarantee accesses to become visible in any particular order. When rendering primitives, quads are guaranteed to be Z-tested, blended and written back in API order (as discussed in detail in [part 9 of this series](https://fgiesen.wordpress.com/2011/07/12/a-trip-through-the-graphics-pipeline-2011-part-9/) (<https://fgiesen.wordpress.com/2011/07/12/a-trip-through-the-graphics-pipeline-2011-part-9/>)), or at least produce the same results as if they were – which takes substantial effort. UAVs make no such effort – UAV accesses happen immediately as they're encountered in the shader, which may be very different from API order. They're not *completely* unordered, though; while there's no guaranteed order of operations within an API call, the API and driver will still collaborate to make sure that perceived sequential ordering is preserved across API calls. Thus, if

you have a complex Compute Shader (or Pixel Shader) writing to an UAV immediately followed by a second (simpler) CS that reads from the same underlying resource, the second CS will see the finished results, never some partially-written output. UAVs support random access. A Pixel Shader can only write to one location per render target – its corresponding pixel. The same Pixel Shader can write to arbitrary locations in whatever UAVs it has bound. UAVs support atomic operations. In the classic Pixel Pipeline, there's no need; we guarantee there's never any collisions anyway. But with the free-form execution provided by UAVs, different threads might be trying to access a piece of memory at the same time, and we need synchronization mechanisms to deal with this.

So from a “CPU programmer”’s point of view, UAVs correspond to regular RAM in a shared-memory multiprocessing system; they’re windows into memory. More interesting is the issue of atomic operations; this is one area where current GPUs diverge considerably from CPU designs.

Atomics

In current CPUs, most of the magic for shared memory processing is handled by the memory hierarchy (i.e. caches). To write to a piece of memory, the active core must first assert exclusive ownership of the corresponding cache line. This is accomplished using what’s called a “cache coherency protocol”, usually MESI (http://en.wikipedia.org/wiki/MESI_protocol) and descendants. The details are tangential to this article; what matters is that because writing to memory entails acquiring exclusive ownership, there’s never a risk of two cores simultaneously trying to write to the same location. In such a model, atomic operations can be implemented by holding exclusive ownership for the duration of the operation; if we had exclusive ownership for the whole time, there’s no chance that someone else was trying to write to the same location while we were performing the atomic operation. Again, the actual details of this get hairy pretty fast (especially as soon as things like paging, interrupts and exceptions get involved), but the 30000-feet-view will suffice for the purposes of this article.

In this type of model, atomic operations are performed using the regular Core ALUs and load/store units, and most of the “interesting” work happens in the caches. The advantage is that atomic operations are (more or less) regular memory accesses, albeit with some extra requirements. There’s a couple of problems, though: most importantly, the standard implementation of cache coherency, “snooping”, requires that all agents in the protocol talk to each other, which has serious scalability issues. There are ways around this restriction (mainly using so-called Directory-based Coherency protocols), but they add additional complexity and latency to memory accesses. Another issue is that all locks and memory transactions really happen at the cache line level; if two unrelated but frequently-updated variables share the same cache line, it can end up “ping-ponging” between multiple cores, causing tons of coherency transactions (and associated slowdown). This problem is called “false sharing”. Software can avoid it by making sure unrelated fields don’t fall into the same cache line; but on GPUs, neither the cache line size nor the memory layout during execution is known or controlled by the application, so this problem would be more serious.

Current GPUs avoid this problem by structuring their memory hierarchy differently. Instead of handling atomic operations inside the shader units (which again raises the “who owns which memory” issue), there’s dedicated atomic units that directly talk to a shared lowest-level cache hierarchy. There’s only one such cache, so the issue of coherency doesn’t come up; either the cache line is present in the cache (which means it’s current) or it isn’t (which means the copy in memory is current). Atomic operations consist of first bringing the respective memory location into the cache (if it isn’t there already), then performing the required read-modify-write operation directly on the cache contents using a dedicated integer ALU on the atomic units. While an atomic unit is busy on a memory location, all other accesses to that location will stall. Since there’s multiple atomic units, it’s necessary to make sure they never try to access the same memory location at the same time; one easy way to accomplish this is to make each atomic unit “own” a certain set of addresses (statically – not dynamically as with cache line ownership). This is done by computing the index of the responsible atomic unit as some hash function of the memory address to be accessed. (Note that I can’t confirm this is how current GPUs do; I’ve found little detail on how the atomic units work in official docs).

If a shader unit wants to perform an atomic operation to a given memory address, it first needs to determine which atomic unit is responsible, wait until it is ready to accept new commands, and then submit the operation (and potentially wait until it is finished if the result of the atomic operation is required). The atomic unit might only be processing one command at a time, or it might have a small FIFO of outstanding requests; and of course there’s all kinds of allocation and queuing details to get right so that atomic operation processing is reasonably fair so that shader units will always make progress. Again, I won’t go into further detail here.

One final remark is that, of course, outstanding atomic operations count as “device memory” accesses, same as memory/textured reads and UAV writes; shader units need to keep track of their outstanding atomic operations and make sure they’re finished when they hit device memory access barriers.

Structured buffers and append/consume buffers

Unless I missed something, these two buffer types are the last CS-related features I haven't talked about yet. And, well, from a hardware perspective, there's not that much to talk about, really. Structured buffers are more of a hint to the driver-internal shader compiler than anything else; they give the driver some hint as to how they're going to be used – namely, they consist of elements with a fixed stride that are likely going to be accessed together – but they still compile down to regular memory accesses in the end. The structured buffer part may bias the driver's decision of their position and layout in memory, but it does not add any fundamentally new functionality to the model.

Append/consume buffers are similar; they could be implemented using the existing atomic instructions. In fact, they kind of are, except the append/consume pointers aren't at an explicit location in the resource, they're side-band data outside the resource that are accessed using special atomic instructions. (And similarly to structured buffers, the fact that their usage is declared as append/consume buffer allows the driver to pick their location in memory appropriately).

Wrap-up.

And... that's it. No more previews for the next part, this series is done :), though that doesn't mean I'm done with it. I have some restructuring and partial rewriting to do – these blog posts are raw and unproofed, and I intend to go over them and turn it into a single document. In the meantime, I'll be writing about other stuff here. I'll try to incorporate the feedback I got so far – if there's any other questions, corrections or comments, now's the time to tell me! I don't want to nail down the ETA for the final cleaned-up version of this series, but I'll try to get it down well before the end of the year. We'll see. Until then, thanks for reading!

From → Coding, Graphics Pipeline

14 Comments

1. Joerg permalink

Very fine granular (per-thread) scheduling is briefly mentioned for the T600 <http://blogs.arm.com/multimedia/534-memory-management-on-embedded-graphics-processors>.

[Reply](#)

2. Alex permalink

Thanks for a fantastic series .. some of the best technical writing around, for its depth, clarity and content.

[Reply](#)

3. wuyuwen permalink

Terrific series, Thanks!

[Reply](#)

4. luke permalink

Thanks for this series. It was a pleasure to read.

[Reply](#)

5. Marcel Ancel permalink

Hi. Very interesting articles. Thanks a lot
and thanks for your demos, great job :)

[Reply](#)

6. jt permalink

Thanks for the series. Do you have a higher level document for the not-so-advanced audience? How does the game content (models, textures) get into the game level? When we launch a game's exe, what goes on in the background? Where do the DX APIs and gpu drivers come into picture when running a game?

I've looked all over the internet but haven't found anything that describes these in an easy to understand way.

[Reply](#)

7. yoelshoshan permalink

Thanks a lot for this great series!
I found it very useful :)

[Reply](#)

8. Augus1990 permalink

Fantastic, thank you for all this info. It would be great if you make a PDF with all the information that you write about Graphics Pipeline, but it's just a suggestion.

I sorry for my english, bye.

Reply

9. royalestel permalink

Holy crap. I thought I was somewhat knowledgeable about computer graphics, since It's been my hobby for a long while now. I see now that I am nowhere close! Thanks for this great series! Cheers!

Reply

10. Tommy permalink

This stuff is gold!

Hope for dx12 vulkan update!

Reply

11. Martijn permalink

A very interesting series. Even though it's currently jan. 2016 it was still very useful for me to read and I learned a lot from it.

I have one question (though I understand if you're done with this series by now, it's been 5 years after all...) about thread groups:

Do I understand correctly that one *and only one* thread group will run on a shader unit at the same time? Because different thread groups would mean they would stomp on each others TGSM I'd think. That would mean you should never make your thread groups too small, or you will get underutilisation of your shader units because of waiting for memory latency, but also not too big or you'll force too many warps onto a single shader unit (and possibly leave whole units unused).

Do I understand this correctly?

I assume this is the same as the local workgroup size in OpenCL, where the documentation gives some vague handwaving like 'the ideal local workgroup size depends on your problem and the hardware' bot no actual insights in how to decide a good group size.

If reading this series gave me the correct insight in this particular point, only that would have been an afternoon well spent (but I learned tons of things from the rest as well).

Reply

o fgiesen permalink

I would've updated it, but GPU architecture (or at least the parts I'm describing) hasn't substantially changed in the last 5 years at all. Microarchitecture yes (the usual incremental improvements), but the big picture has been pretty stable.

Anyway, "shader units" (reading this to mean EUs for Intel GPUs, CUs for AMD GCN and SMXs for NVidia GPUs) can run multiple thread groups at the same time no problem. Thread group memory is allocated in a similar way to GPU registers: the hardware units have a fixed-size pool, and the maximum number of active thread groups on a shader unit depends (among other things) on how many of them fit in said pool at the same time. (Subject to a certain allocation granularity which depends on the hardware, and also subject to other restrictions that again depend on the hardware).

This series purposefully stays away from shader core internals since that was very much in flux at the time of writing, and besides, it's the part that the vendors tend to tell you lots of details about anyway.

In general, you want to keep your thread groups relatively small if possible. Around 64 threads is a good value for current AMD and NV hardware (Intel prefers less). AMD GCN hardware cannot run less than 64 threads (1 wavefront) in a thread group; make them any smaller and they'll still allocate 64 threads worth of resources to it. Current NV hardware has a similar limit (warp size), though theirs is 32 threads (though 64 is still often better). Other than that, a thread group is a unit of threads that can communicate efficiently. Depending on the task, sometimes you want more. If you don't care much, staying close to 64 is a good rule of thumb at the moment.

Reply

12. Chris permalink

It tooks me three hours to find I was looking for. Thank you very much for this interesting and profound article about GPU parallelity!

Reply

Trackbacks & Pingbacks

1. A trip through the Graphics Pipeline 2011: Index « The ryg blog

Blog at WordPress.com.